

06-01-00

A

**UTILITY PATENT APPLICATION TRANSMITTAL**  
**(Large Entity)***(Only for new nonprovisional applications under 37 CFR 1.53(b))*Docket No.  
POU9-2000-0014-US1

Total Pages in this Submission

**TO THE ASSISTANT COMMISSIONER FOR PATENTS**Box Patent Application  
Washington, D.C. 20231

Transmitted herewith for filing under 35 U.S.C. 111(a) and 37 C.F.R. 1.53(b) is a new utility patent application for invention entitled:

**METHOD, SYSTEM AND PROGRAM PRODUCTS FOR SERIALIZING REPLICATED  
TRANSACTIONS OF A DISTRIBUTED COMPUTING ENVIRONMENT**

and invented by:

**Marcos N. Novaes, Gregory D. Laib, Jeffrey S. Lucash, Rosario A. Uceda-Sosa**If a **CONTINUATION APPLICATION**, check appropriate box and supply the requisite information:☐ Continuation ☐ Divisional ☐ Continuation-in-part (CIP) of prior application No.: \_\_\_\_\_

Which is a:

☐ Continuation ☐ Divisional ☐ Continuation-in-part (CIP) of prior application No.: \_\_\_\_\_

Which is a:

☐ Continuation ☐ Divisional ☐ Continuation-in-part (CIP) of prior application No.: \_\_\_\_\_

Enclosed are:

**Application Elements**

1. ☒ Filing fee as calculated and transmitted as described below
2. ☒ Specification having 41 pages and including the following:
  - a. ☒ Descriptive Title of the Invention
  - b. ☒ Cross References to Related Applications *(if applicable)*
  - c. ☐ Statement Regarding Federally-sponsored Research/Development *(if applicable)*
  - d. ☐ Reference to Microfiche Appendix *(if applicable)*
  - e. ☒ Background of the Invention
  - f. ☒ Brief Summary of the Invention
  - g. ☒ Brief Description of the Drawings *(if drawings filed)*
  - h. ☒ Detailed Description
  - i. ☒ Claim(s) as Classified Below
  - j. ☒ Abstract of the Disclosure

# UTILITY PATENT APPLICATION TRANSMITTAL (Large Entity)

(Only for new nonprovisional applications under 37 CFR 1.53(b))

Docket No.  
POU9-2000-0014-US1

Total Pages in this Submission

## Application Elements (Continued)

3. ☒ Drawing(s) (when necessary as prescribed by 35 USC 113)

a. ☒ Formal Number of Sheets 22

b. ☐ Informal Number of Sheets \_\_\_\_\_

4. ☒ Oath or Declaration

a. ☐ Newly executed (original or copy) ☒ Unexecuted

b. ☐ Copy from a prior application (37 CFR 1.63(d)) (for continuation/divisional application only)

c. ☒ With Power of Attorney ☐ Without Power of Attorney

d. ☐ DELETION OF INVENTOR(S)

Signed statement attached deleting inventor(s) named in the prior application,  
see 37 C.F.R. 1.63(d)(2) and 1.33(b).

5. ☐ Incorporation By Reference (usable if Box 4b is checked)

The entire disclosure of the prior application, from which a copy of the oath or declaration is supplied  
under Box 4b, is considered as being part of the disclosure of the accompanying application and is hereby  
incorporated by reference therein.

6. ☐ Computer Program in Microfiche (Appendix)

7. ☐ Nucleotide and/or Amino Acid Sequence Submission (if applicable, all must be included)

a. ☐ Paper Copy

b. ☐ Computer Readable Copy (identical to computer copy)

c. ☐ Statement Verifying Identical Paper and Computer Readable Copy

## Accompanying Application Parts

8. ☐ Assignment Papers (cover sheet & document(s))

9. ☐ 37 CFR 3.73(B) Statement (when there is an assignee)

10. ☐ English Translation Document (if applicable)

11. ☒ Information Disclosure Statement/PTO-1449 ☒ Copies of IDS Citations

12. ☐ Preliminary Amendment

13. ☒ Acknowledgment postcard

14. ☒ Certificate of Mailing

☐ First Class ☒ Express Mail (Specify Label No.): EL643175337US

**UTILITY PATENT APPLICATION TRANSMITTAL**  
**(Large Entity)**

*(Only for new nonprovisional applications under 37 CFR 1.53(b))*

Docket No.  
**POU9-2000-0014-US1**

Total Pages in this Submission

**Accompanying Application Parts (Continued)**

15. ☐ Certified Copy of Priority Document(s) *(if foreign priority is claimed)*

16. ☐ Additional Enclosures *(please identify below):*

**Fee Calculation and Transmittal**

**CLAIMS AS FILED**

For	#Filed	#Allowed	#Extra	Rate	Fee
Total Claims	3	- 20 =	0	x \$18.00	\$0.00
Indep. Claims	3	- 3 =	0	x \$78.00	\$0.00
Multiple Dependent Claims (check if applicable) <input type="checkbox"/>					\$0.00
BASIC FEE					\$690.00
OTHER FEE (specify purpose)					\$0.00
TOTAL FILING FEE					\$690.00

- ☐ A check in the amount of \_\_\_\_\_ to cover the filing fee is enclosed.
- ☒ The Commissioner is hereby authorized to charge and credit Deposit Account No. **09-0463 (IBM)** as described below. A duplicate copy of this sheet is enclosed.
- ☒ Charge the amount of **\$690.00** as filing fee.
  - ☒ Credit any overpayment.
  - ☒ Charge any additional filing fees required under 37 C.F.R. 1.16 and 1.17.
  - ☐ Charge the issue fee set in 37 C.F.R. 1.18 at the mailing of the Notice of Allowance, pursuant to 37 C.F.R. 1.311(b).

Blanche E. Schiller  
Signature

Dated: **May 30, 2000**

**Blanche E. Schiller, Esq.**  
**Reg. No. 35,670**  
**HESLIN & ROTHENBERG, P.C.**  
**5 Columbia Circle**  
**Albany, NY 12203**  
**Telephone: (518) 452-5600**  
**Facsimile: (518) 452-5579**

cc:

	1990	1991	1992	1993	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032	2033	2034	2035	2036	2037	2038	2039	2040	2041	2042	2043	2044	2045	2046	2047	2048	2049	2050	2051	2052	2053	2054	2055	2056	2057	2058	2059	2060	2061	2062	2063	2064	2065	2066	2067	2068	2069	2070	2071	2072	2073	2074	2075	2076	2077	2078	2079	2080	2081	2082	2083	2084	2085	2086	2087	2088	2089	2090	2091	2092	2093	2094	2095	2096	2097	2098	2099	2100	2101	2102	2103	2104	2105	2106	2107	2108	2109	2110	2111	2112	2113	2114	2115	2116	2117	2118	2119	2120	2121	2122	2123	2124	2125	2126	2127	2128	2129	2130	2131	2132	2133	2134	2135	2136	2137	2138	2139	2140	2141	2142	2143	2144	2145	2146	2147	2148	2149	2150	2151	2152	2153	2154	2155	2156	2157	2158	2159	2160	2161	2162	2163	2164	2165	2166	2167	2168	2169	2170	2171	2172	2173	2174	2175	2176	2177	2178	2179	2180	2181	2182	2183	2184	2185	2186	2187	2188	2189	2190	2191	2192	2193	2194	2195	2196	2197	2198	2199	2200	2201	2202	2203	2204	2205	2206	2207	2208	2209	2210	2211	2212	2213	2214	2215	2216	2217	2218	2219	2220	2221	2222	2223	2224	2225	2226	2227	2228	2229	2230	2231	2232	2233	2234	2235	2236	2237	2238	2239	2240	2241	2242	2243	2244	2245	2246	2247	2248	2249	2250	2251	2252	2253	2254	2255	2256	2257	2258	2259	2260	2261	2262	2263	2264	2265	2266	2267	2268	2269	2270	2271	2272	2273	2274	2275	2276	2277	2278	2279	2280	2281	2282	2283	2284	2285	2286	2287	2288	2289	2290	2291	2292	2293	2294	2295	2296	2297	2298	2299	2300	2301	2302	2303	2304	2305	2306	2307	2308	2309	2310	2311	2312	2313	2314	2315	2316	2317	2318	2319	2320	2321	2322	2323	2324	2325	2326	2327	2328	2329	2330	2331	2332	2333	2334	2335	2336	2337	2338	2339	2340	2341	2342	2343	2344	2345	2346	2347	2348	2349	2350	2351	2352	2353	2354	2355	2356	2357	2358	2359	2360	2361	2362	2363	2364	2365	2366	2367	2368	2369	2370	2371	2372	2373	2374	2375	2376	2377	2378	2379	2380	2381	2382	2383	2384	2385	2386	2387	2388	2389	2390	2391	2392	2393	2394	2395	2396	2397	2398	2399	2400	2401	2402	2403	2404	2405	2406	2407	2408	2409	2410	2411	2412	2413	2414	2415	2416	2417	2418	2419	2420	2421	2422	2423	2424	2425	2426	2427	2428	2429	2430	2431	2432	2433	2434	2435	2436	2437	2438	2439	2440	2441	2442	2
--	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	---

Title: METHOD, SYSTEM AND PROGRAM PRODUCTS FOR SERIALIZING  
REPLICATED TRANSACTIONS OF A DISTRIBUTED COMPUTING  
ENVIRONMENT

Date of Deposit May 31, 2000

BOX PATENT APPLICATION  
ASSISTANT COMMISSIONER FOR PATENTS  
WASHINGTON, D.C. 20231

(Typed or printed name of person mailing paper or fee)

Jill K. Becker  
person mailing paper on

(Signature of ~~person~~ mailing paper or fee)

- \* U.S. Patent Application which includes:
  - Specification (37 pages), 3 Claims (3 pages), Abstract (1 page), and twenty-two (22) sheets of Formal Drawings
- \* Utility Patent Application Transmittal Letter (3 pages) (in duplicate)
- \* Declaration and Power of Attorney for Patent Application (4 pages) (Unsigned)
- \* Information Disclosure Citation (1 page) and six (6) references
- \* Two (2) Acknowledgment Postcards

**METHOD, SYSTEM AND PROGRAM PRODUCTS FOR  
SERIALIZING REPLICATED TRANSACTIONS OF A  
DISTRIBUTED COMPUTING ENVIRONMENT**

**Cross-Reference to Related Applications**

5        This application contains subject matter which is  
related to the subject matter of the following applications,  
each of which is assigned to the same assignee as this  
application and filed on the same day as this application.  
Each of the below listed applications is hereby incorporated  
10    herein by reference in its entirety:

      "METHOD, SYSTEM AND PROGRAM PRODUCTS FOR MANAGING  
PROCESSING GROUPS OF A DISTRIBUTED COMPUTING ENVIRONMENT,"  
Novaes et al., (Docket No. POU9-2000-0003-US1), Serial No.  
\_\_\_\_\_, filed herewith;

15        "METHOD, SYSTEM AND PROGRAM PRODUCTS FOR RECOVERING  
FROM FAILURES WITHIN A SHARED NOTHING DISTRIBUTED COMPUTING  
ENVIRONMENT," Novaes et al., (Docket No. POU9-2000-0009-  
US1), Serial No. \_\_\_\_\_, filed herewith;

20        "SYNCHRONOUS REPLICATION OF TRANSACTIONS IN A  
DISTRIBUTED SYSTEM," Novaes et al., (Docket No. POU9-2000-  
0006-US1), Serial No. \_\_\_\_\_, filed herewith; and

"METHOD, SYSTEM AND PROGRAM PRODUCTS FOR MANAGING A CLUSTERED COMPUTING ENVIRONMENT," Novaes et al., (Docket No. POU9-2000-0004-US1), Serial No. \_\_\_\_\_, filed\_\_\_\_\_.

### **Technical Field**

5           This invention relates, in general, to distributed systems, and in particular, to managing a distributed synchronous transaction system.

### **Background Art**

10           Distributed systems are highly-available, scalable systems that are utilized in various situations, including those situations that require a high-throughput of work or continuous or nearly continuous availability of the system.

15           One type of a distributed system is a distributed synchronous transaction system, which is a system that performs distributed synchronous transactions on behalf of distributed clients. A distributed synchronous transaction is a transaction that is initiated substantially immediately when it is requested by a client application, and which in turn, is notified of the success of the transaction  
20           substantially immediately following the transaction's completion.

          Although there are facilities today for managing distributed synchronous transactions, these facilities tend to be complicated. Thus, there is still a need for

capabilities to facilitate the management of synchronous transactions in a distributed system.

### Summary of the Invention

The shortcomings of the prior art are overcome and additional advantages are provided through the provision of a method of serializing replicated transactions in a distributed computing environment. The method includes, for instance, initiating a modification operation on a resource of a distributed computing environment; during a phase of the modification operation, detecting whether a conflict for the resource exists; and satisfying the conflict, if the conflict exists, without requiring explicit locking of the resource.

System and computer program products corresponding to the above-summarized methods are also described and claimed herein.

Additional features and advantages are realized through the techniques of the present invention. Other embodiments and aspects of the invention are described in detail herein and are considered a part of the claimed invention.

### **Brief Description of the Drawings**

The subject matter which is regarded as the invention is particularly pointed out and distinctly claimed in the claims at the conclusion of the specification. The  
5 foregoing and other objects, features, and advantages of the invention are apparent from the following detailed description taken in conjunction with the accompanying drawings in which:

FIG. 1 depicts one example of a computing  
10 environment incorporating and using aspects of the present invention;

FIG. 2 depicts one example of various components of several nodes of FIG. 1, in accordance with an aspect of the present invention;

FIG. 3 depicts one embodiment of a computing  
15 environment in which a client application instance replies to a request of a third party application without using a DSTS server, in accordance with an aspect of the present invention;

FIG. 4 depicts one embodiment of a computing  
20 environment in which a client application instance uses a DSTS server to reply to a request of the third party application, in accordance with an aspect of the present invention;

FIG. 5 depicts one example of processing group, used in accordance with an aspect of the present invention;

Fig. 6a depicts one example of the components associated with a group activation protocol, in accordance with an aspect of the present invention;

FIGs. 6b-6d depict one embodiment of the logic associated with performing group activation, in accordance with an aspect of the present invention;

FIG. 7 depicts one example of the fields associated with an initialize message, in accordance with an aspect of the present invention;

FIG. 8 depicts one embodiment of the components associated with a group join protocol, in accordance with an aspect of the present invention;

FIGs. 9a-9b depict one embodiment of the logic associated with joining a processing group, in accordance with an aspect of the present invention;

FIG. 10 depicts one example of the fields associated with a quiesce message, in accordance with an aspect of the present invention;

FIG. 11 depicts one embodiment of the fields associated with an archive message, in accordance with an aspect of the present invention;

5 FIG. 12 depicts one embodiment of the fields associated with a dearchive message, in accordance with an aspect of the present invention;

FIG. 13 includes one example of the fields associated with an enumerate handles message, in accordance with an aspect of the present invention;

10 FIG. 14 depicts one example of the fields associated with a handle enumeration message, in accordance with an aspect of the present invention;

15 FIG. 15 depicts one embodiment of the logic associated with excluding a member from a processing group, in accordance with an aspect of the present invention;

FIG. 16 depicts one example of the fields associated with a quorum notification message, in accordance with an aspect of the present invention;

20 FIG. 17 depicts one example of the fields associated with a replicate request message, in accordance with an aspect of the present invention;

FIG. 18 depicts one example of the fields associated with a replication callback message, in accordance with an aspect of the present invention;

FIG. 19 depicts one example of the fields associated with a replication callback result message, in accordance with an aspect of the present invention;

FIG. 20 depicts one example of the fields associated with a replication completed message, in accordance with an aspect of the present invention;

FIG. 21 depicts one example of the fields associated with a shutdown message, in accordance with an aspect of the present invention;

FIGs. 22a-22b depict one embodiment of the flow of messages associated with processing a synchronous transaction, in accordance with an aspect of the present invention;

FIG. 23 depicts one embodiment of the flow of messages associated with a Prepare to Commit operation, in accordance with an aspect of the present invention;

FIG. 24 depicts one embodiment of the message flow associated with a Commit operation, in accordance with an aspect of the present invention;



operating system. Each processing node within a frame is coupled to the other processing nodes of the frame via, for example, an internal LAN connection. Additionally, each frame is coupled to the other frames via LAN gates 104.

5       As examples, each LAN gate 104 includes either a RISC/6000 computer, any computer network connection to the LAN, or a network router. However, these are only examples. It will be apparent to those skilled in the relevant art that there are other types of LAN gates, and that other  
10 mechanisms can also be used to couple the frames to one another.

The distributed computing environment of FIG. 1 is only one example. It is possible to have more or less than eight frames, or more or less than sixteen nodes per frame.

15 Further, the processing nodes do not have to be RISC/6000 computers running AIX. Some or all of the processing nodes can include different types of computers and/or different operating systems. Further, a heterogeneous environment can include and utilize the invention, in which one or more of  
20 the nodes and/or operating systems of the environment are distinct from other nodes or operating systems of the environment. The nodes of such a heterogeneous environment interoperate, in that they collaborate and share resources with each other, as described herein.

25       Further details regarding the nodes of a distributed computing environment are described with reference to FIG. 2. In one example, a distributed client application 200

runs on a plurality of nodes 202. In particular, an instance of the client application runs substantially simultaneously on each of the plurality of nodes, which includes three nodes in this specific example. (It will be  
5 apparent to one skilled in the art that the client application can run on any number of the nodes of the environment, including only one node.)

In one embodiment, the client application instances are coupled to a distributed synchronous transaction system  
10 (DSTS), which enables the application instances, in accordance with an aspect of the present invention, to participate in the synchronous replication of transactions. By using the distributed synchronous transaction system, a client instance is able to participate in synchronous  
15 replication of transactions, even though the client application instance has no direct knowledge of any other instances of the application. The distributed synchronous transaction system includes one or more DSTS instances (e.g., computer programs) 204 that run on one or more nodes.  
20 In one example, a DSTS instance is executed on each node that has a client application instance interested in participating in a distributed transaction. Each DSTS instance is coupled to one or more instances of one or more client applications.

25 When the DSTS instance is loaded into a node's memory and executed, it is perceived as a server process, which serves its corresponding client application process (or processes). It is the DSTS system that performs a

distributed synchronous transaction on behalf of a client application. When the transaction is requested by the client, it is initiated substantially immediately by a DSTS server. Further, the client is substantially immediately  
5 notified of the outcome (e.g., success, failure) of the transaction, upon completion of the transaction.

A collection of one or more client application instances participating in the execution of a distributed synchronous transaction is referred to as a replicated group  
10 of client application instances. This group is distinct from other forms of groups in a distributed system, since the members of the replicated group have no direct knowledge of one another. Instead, the group is implicitly formed, when a client application instance diverts a flow of update  
15 operations to be replicated to one or more other client application instances.

In particular, the client application diverts the flow of operations, which modify its persistent (stored) or run-time (not stored) state. These update operations are  
20 classified as write operations. Any other transaction which does not modify the state of the client application can be termed a query, or read transaction. In accordance with an aspect of the present invention, client applications perform write operations as distributed synchronous transactions,  
25 which provides each copy of the client application with a consistent, or identical state. Such capability in turn makes possible for any copy of the application to respond to queries (read operations) to its state without having to

redirect the query to any of the other replicas. In other words, client applications can service read operations locally without using a DSTS server (see FIG. 3), while write operations are replicated to other instances of the client application, and thus, use DSTS (see FIG. 4), as described in further detail below. This architecture is optimal for, but not limited to, systems which are read intensive, and that exhibit a low rate of write operations.

The flow of update operations is diverted by a client application via, for instance, a DSTS protocol used by the client application. One feature of this protocol, in accordance with an aspect of the present invention, includes membership in one or more processing groups. A processing group 500 (FIG. 5) includes one or more members 502. Each member, in this example, is a DSTS server. Thus, for each client application instance of a replicated group, there is a corresponding DSTS server in a given processing group (a.k.a., a group). For example, if a replicated group includes Client Application Instances A and B, then a processing group includes DSTS Servers A and B, which are coupled to Application Instances A and B, respectively. This allows the processing group to handle the replication of transactions for the client applications of the replicated group, and enables the replication to be transparent to those client applications.

Each member of a processing group is ensured a consistent view of the group's state data. The data is kept consistent because it is only updated by well-defined group

protocols. Examples of the protocols include admission to a group, including activation of the group and joining the group, and exclusion from the group, each of which is described in detail below. Further details regarding the management of a processing group are discussed in U.S. Patent No. 5,748,958 entitled "System For Utilizing Batch Requests To Present Membership Changes To Process Groups," issued on May 5, 1998, which is hereby incorporated herein by reference in its entirety.

One embodiment of the logic associated with admission to a group is described with reference to FIGs. 6a-6d. In particular, FIG. 6a depicts one example of the components involved in activating a group; and FIGs. 6b-6d depict one embodiment of the logic. In the initial case of group activation, there are no members in the processing group. The group is assumed to have been previously defined, but none of the copies (i.e., DSTS) of the group are currently being executed. A DSTS copy begins to be executed, when it is connected to by a client application.

In one example, a client application 602 connects to a DSTS server 604 via an initialize message, STEP 600 (FIGs. 6a, 6b). The initialize message is sent from client application instance 602 to DSTS server 604 to connect to the DSTS system. Specifically, in one example, the client application instance connects to the DSTS server on the same node as the client application instance. One example of the initialize message is described with reference to FIG. 7.

An initialize message 700 includes, for instance, an operation code 702 indicating the type of operation (e.g., initialize) being requested, and a name 704 of the client application issuing the request. The DSTS system uses the application name to propagate transactions to the other instances of the application (i.e., the members of the replicated group) having the same name.

Referring back to FIGs. 6a-6b, in response to this message, the DSTS server proposes to join a group (designed by application name 704 (FIG. 7), STEP 606 (FIG. 6b). As it proposes to join the group, the DSTS server reads the group state from persistent storage 608 (FIG. 6a). The group state 610 includes, for instance, the group sequence number and the activation status. If the group state is active, INQUIRY 612 (FIG. 6b), the joining copy executes a join protocol, STEP 614, as described below. Otherwise, the state is inactive, and the copy is able to join the group immediately, without executing the below defined join protocol, STEP 616.

As the DSTS server joins the group, the copy compares the group's sequence number with its own sequence number, STEP 618. If the group's sequence number is smaller than its own, then the copy updates the group's sequence number, STEP 620. Thereafter, or if the group's sequence number is equal to or larger than the copy's sequence number, a determination is made as to whether a quorum (in this example) of members has been reached, INQUIRY 622.

If quorum has not been reached, then processing continues with STEP 600, for another member, at least until quorum is reached. As a quorum of members join the group, the copies which are members of the processing group have  
5 knowledge that the quorum was achieved. At this point, the group's sequence number is set to the highest incarnation of the members, STEP 624. The members, whose sequence number match the group's when this point is reached, initiate an activation protocol by sending a group activation message,  
10 STEP 626. The group activation message initiates a multi-phase protocol.

In the first phase of activation, the members of the group receive the group activation message, which contains the node address of the member which sent the message, STEP  
15 628 (FIG. 6c). Then, the current group members whose sequence numbers are lower than the current group's sequence number ask the sender of the activation message for a copy of the group state that is associated with the group's sequence number, STEP 630. These members reinitialize  
20 themselves using the new group state, STEP 632, and then propose to continue to the second phase of group activation, STEP 634. Any member that fails initialization at this point votes to abort the protocol.

The members whose sequence number match those of the  
25 group also propose to go to the second phase. If all current members propose to go to the second phase (none aborts), the second phase begins.

As the first phase of group activation finishes, the current members of the processing group verify that a majority of the members was maintained, STEP 636 (FIG. 6d). Furthermore, each member now has the same consistent  
5 sequence number and copy of the distributed state.

The members now change the group sequence number by, for instance, adding 1 to it, STEP 638. The members then store the new sequence number in group state and propose to conclude the protocol, STEP 640. Any member that fails at  
10 this stage proposes to abort the protocol.

In protocol completion, if no current member aborted, INQUIRY 642, then the group has the guarantee that the current members of the group have the same consistent group state and sequence number, and that the new sequence number  
15 has been stored by a majority of the numbers of the group. The group state is then changed to active, STEP 644.

Each time a member joins an active group, it initiates a multi-phase group admission protocol, one embodiment of which is described with reference to FIGs. 8 and 9a-9b. In  
20 particular, FIG. 8 depicts the components of the join process, while FIGs. 9a-9b depict one embodiment of the logic. In the first phase of the protocol, the joining member (800 of FIG. 8) sends a join proposal message with the sequence number that it retrieved from persistent  
25 storage, or a negative infinity, if it was unable to retrieve the sequence number, STEP 900 (FIG. 9a). As examples, the sequence number, as well as other group state,

may not be available, when the disk where the state is stored is corrupted or is otherwise not available, or when this is actually the first time that the member copy is being executed under any given processor.

5        In response to receiving the join proposal message, the other members of the group (802, FIG 8) cease to make updates to the distributed data, STEP 902. In one embodiment, in order to cease the updates, each member of the group sends a quiesce message to its corresponding  
10 client application instance. One example of the quiesce message is described with reference to FIG. 10.

A quiesce message 1000 includes, for instance, an operation code 1002 specifying that this a quiesce operation. The quiesce message requests the client  
15 applications to cease sending update requests (e.g., replicate request messages described below), such that the global state of the application is stabilized.

Thereafter, each copy of the application is requested to produce a snapshot of the current state of the  
20 application and to store this state in persistent storage, STEP 904. This request is performed by sending an archive message to the copies of the application. One example of an archive message is described with reference to FIG. 11. In one example, an archive message 1100 includes an operation  
25 code 1102 indicating that this is an archive request.



An enumerate handles message 1300 (FIG. 13) includes, for example, an operation code 1302 indicating that this is an enumerate handles message. After receiving this message, the client application returns a handle enumeration message  
5 to the DSTS system, which maps the names of the resources that the application has created to resource handles.

One example of the handle enumeration message is described with reference to FIG. 14 and includes, for example, an operation code 1402 indicating that this is the  
10 handle enumeration message, and a resource handle map 1404, which includes one or more pairs of resource names and handles. These handles are unique names used, for instance, to notify third party applications of changes to the client application's state, and to serialize simultaneous update  
15 requests to the same resources, as described below.

After successfully reinitializing itself by loading the snapshot, the new copy is allowed to participate in the DSTS system, and a resume message is sent to all copies such that the DSTS system may resume normal operation. Further, the  
20 new copy proposes to begin the second phase of join, STEP 912.

Returning to INQUIRY 908, if the group becomes inactive, the joining member notes the fact that its sequence number is outdated, STEP 916, and waits for an  
25 activation message to take further action, STEP 918. The joining member does not take place in the second phase of join.

Returning to INQUIRY 906, if the joining member's sequence number is equal to the sequence number of the group, then the group is inactive. This fact is given by a virtue of the group activation protocol (e.g., a quorum  
5 policy, in this example) and by the property of quorum enforcement. Thus, the joining member waits for an activation message to take effect, STEP 918, and there is no second phase of join. Similarly, if the joining member's sequence number is higher, INQUIRY 906, it also follows that  
10 the group is inactive, and thus, the joining member waits for an activation message, STEP 918.

If the joining member has proposed to proceed to the second phase, it has the new sequence number and distributed state. Thus, the members (including the joining member) now  
15 change the group's sequence number by, for instance, adding one to it, STEP 922 (FIG. 9b). The members then store the new sequence number and group state, STEP 924, and further, they propose to conclude the protocol, STEP 926. Any member that fails at this stage, proposes to abort the protocol.  
20 If no member aborts, the group is guaranteed that the current members of the group have the same consistent group state and sequence number, and that the new sequence number has been stored for a majority of the members of the group.

In addition to the above, a member can be excluded from  
25 a group. In particular, each time a node fails, or the DSTS copy that executes on the node fails, the remaining members of the group notice that a member has failed, STEP 1500 (FIG. 15). If the group is inactive, INQUIRY 1502, no

action is taken, STEP 1504. Further, if the group is active, but does not have a majority of members, INQUIRY 1506, then no action is taken.

However, if the group is active and retains majority,  
5 INQUIRY 1506, then each member stops any further updates to the distributed state, STEP 1507. Additionally, each member changes the group sequence number by, for instance, adding 1 to it, STEP 1508, and stores the new sequence number and the group state, STEP 1510. Then, the members propose to  
10 conclude the protocol, STEP 1512. Any member that fails at this stage proposes to abort the protocol.

If no member aborts, then the group has a guarantee that the current members of the group have the same consistent group state and sequence number, and that the new  
15 sequence number has been stored by a majority of the members of the group.

The DSTS system notifies the client application instances when a quorum (majority) of DSTS servers is available or has been lost, by utilizing, for instance, a  
20 quorum notification message. In one example, a quorum notification message 1600 (FIG. 16) includes an operation code 1602, and the quorum information 1604, indicating whether the group has quorum.

As described herein, members of a processing group are  
25 utilized to replicate distributed synchronous transactions, which are initiated by client application instances coupled

to the members of the group. To facilitate communication between the client instances and the server members of the group, various messages are employed. In one example, these messages include (in addition to the messages described  
5 above) a replicate request message, a replication callback message, a replication callback result message, a replication completed message and a shutdown message, each of which is described below.

One example of a replicate request message is described  
10 with reference to FIG. 17. A replicate request message 1700 is a message that initiates the distributed transaction. In one example, it includes an operation code 1702 indicating that this is a replicate request message; a list of the new resource names 1704 being created, if any; an exclusive  
15 access set 1706 specifying zero or more exclusive resources of the client application; a shared access set 1708 specifying zero or more shared resources of the client application; a replication policy 1710 providing rules to be adhered to during the replication (e.g., a quorum of the  
20 group needed to proceed with certain tasks); a request 1712 specifying the transaction to be replicated and performed (e.g., a create or update request); and a request size 1714 indicating the size of the request.

The replicate request message is sent by a single  
25 client application instance (a.k.a., the initiator) to a server process of the DSTS system. Upon receipt of the message (or sometime thereafter), the server process distributes the message to one or more other server

processes of the distributed computing environment. In particular, in one example, it is sent to all of the other current server processes of the processing group.

In response, each of the server processes sends a  
5 replication callback message to the corresponding instances (peers) of the client application. One example of a replication callback message is described with reference to FIG. 18. A replication callback message 1800 includes, for instance, an operation code 1802 indicating that this is a  
10 replication callback message; an array of the new resource names 1804, if any are to be created; an exclusive access set 1806 specifying zero or more exclusive resources of the client application; a shared access set 1808 specifying zero or more shared resources of the client application; a  
15 request 1810 specifying the transaction to be replicated and performed; and a request size 1812 indicating the size of the request.

In addition to the above, a replication callback result message is sent from the client application to the DSTS  
20 server, after the requested transaction is processed. One example of a replication callback result messages is described with reference to FIG. 19. A replication callback result message 1900 includes an operation code 1902 indicating that this is a replication callback result  
25 message; an array of the new resource names 1904, if any, along with their handles (e.g., unique identifiers); a modified resource set 1906, including the handles of any

modified resources; and a deleted resource set 1908,  
including the handles of any deleted resources.

After the server processes receive the replication  
callback results, they verify that the transaction has been  
5 completed by forwarding a replication completed message 2000  
(FIG. 20). In one example, replication completed message  
2000 includes an operation code 2002 indicating that this is  
a replication completed message; and an operation status  
2004 specifying whether the transaction was performed  
10 successfully.

Should the system be shut down, the DSTS system  
utilizes a shutdown message that notifies the copies of the  
client application that the system is about to shut down. In  
one example, a shutdown message 2100 (FIG. 21) includes an  
15 operation code 2102 indicating that shutdown is to be  
performed. This message has the objective of allowing the  
copies of the client application to perform a graceful  
shutdown procedure, terminating any pending transaction(s).  
When the client applications terminate the shutdown process,  
20 they reply with a shutdown acknowledgment to the DSTS  
system.

Utilization of the above-described replication messages  
is further described below with reference to FIGs. 22a and  
22b. Referring to FIG. 22a, a replicate request message  
25 2200 is sent by a single client application instance 2202 to  
a server process 2204 of the DSTS system. The server then  
distributes 2206 the replicate request message to the other





cases, the tokens do not conflict, so there is a great improvement in performance over a token granting server approach. But in the case in which tokens do conflict, the serialization technique of the present invention is  
5 performed in order to preserve the consistency of the data in each member of the processing server group.

For example, assume that two transactions are simultaneously initiated, that request exclusive access to a token labeled "A". Further, assume that Server 1 initiates  
10 transaction T1, and Server 2 initiates transaction T2. Assume that T1 is supposed to set A=1 and T2 is to set A=2. Assume further there are three members in the processing group, which are to perform these transactions. Since the transactions are initiated simultaneously, their order is  
15 not important, but they are to be executed in the same order by all the members.

The synchronously replicated transactions are executed using a two-phase commit protocol. Thus, the data is transmitted in a first phase, called the Prepare to Commit  
20 (PTC) phase, and the transaction is committed in a second phase, called the Commit (CMT) phase. The two-phase commit can proceed in parallel (i.e., transactions T1 and T2 can be initiated in parallel), allowing the replication of transactions to be more efficient. However, at some point  
25 in the two-phase commit protocol, the transactions are to be serialized. If not, problems arise, as described below.

If the two-phase commit is allowed to proceed in parallel without serialization, it could lead to inconsistent results, as illustrated below:

	<u>Server 1</u>	<u>Server 2</u>	<u>Server 3</u>
5	PTC(T1)	PTC(T2)	PTC(T2)
	PTC(T2)	PTC(T1)	PTC(T1)

/\*\*the servers wait for acknowledgment that the PTCs were received before processing the Commit phase:

	CMT(T1)	CMT(T1)	CMT(T2)
10	CMT(T2)	CMT(T2)	CMT(T1)

The problem here is that Server 1 and Server 2 executed T1, T2, setting A=1, in these servers. However, Server 3 executed T2, T1, setting A=2, as a final result. The value of "A" is now inconsistent in the processing group, and that is not acceptable in a synchronously replicated transaction system.

In order to overcome this problem, the first phase of the two-phase commit process (the PTC phase) is allowed to proceed in parallel, and then the Commit phase is serialized based on the token information sent in the PTC, in accordance with an aspect of the present invention. The PTC protocol is extended such that it provides information on which tokens are necessary for exclusive/shared access for each transaction. Since an assignment (A=1) requires

exclusive access, the token "A" is listed for exclusive access in the PTC of both T1 and T2.

Further details relating to the two-phase commit protocol is described with reference to FIGs. 23 and 24. In particular, one example of the first phase of the two-phase commit protocol, the Prepare to Commit phase, is described with reference to FIG. 23, and one example of the second phase, the Commit phase, is described with reference to FIG. 24.

Referring to FIG. 23, initially, a replicate request message 2300 is sent from client application instance 2302 to server 2304 indicating that a PTC is to be performed. In response to receiving the PTC request, server 2304 sends a PTC message 2306 to the other servers of the group (e.g., server 2308a and 2308b). In one example, the PTC message includes the same fields as the replicate request message, as well as an identifier of the request. Since server 2304 is initiating the PTC, it is referred to as the protocol initiator.

Thereafter, each non-initiator server responds to the PTC request with a PTC acknowledgment (PTC\_ACK) message 2310. In particular, server 2308a sends an acknowledgment, which includes an operation code, as well as the request identifier. Similarly, server 2308b sends an acknowledgment, but only after serializing any conflicts. That is, in this example, server 2308b is chosen as a coordinator of the group. Thus, it monitors all of the PTC

requests it receives and sends a PTC\_ACK message 2310 serializing any conflicting requests. If it notices that two or more PTCs are issued for the same exclusive access resource (or for an exclusive request which conflicts with a shared one), then the group coordinator chooses to commit one of them first, waits for the confirmation that the update is complete, and then commits the second one, and so forth.

The protocol initiator (e.g., server 2304) receives the PTC\_ACK messages from the other servers. After it receives all of the PTC\_ACK messages for a given message, it sends a commit message, thus, initiating the second phase of the two-phase commit protocol.

One example of the second phase of the two-phase commit protocol is described with reference to FIG. 24. Initially, the protocol initiator 2400 receives PTC\_ACK messages from all of the members of the group, and then sends a commit message 2402 to each of the other servers of the processing group. Each server of the group sends a replication callback message 2404 to its corresponding application to request the application to commit the operation. After committing the operation, a replication callback result message 2406 is sent from the client application to the DSTS server.

Thereafter, a commit acknowledge message 2408 is sent from each DSTS server to the protocol initiator (e.g., server 2400). The protocol initiator receives the commit

acknowledge messages from all the members of the group and sends a replication completed message 2410 to the initiating client, if at least a majority of the members have completed the request.

5           In accordance with an aspect of the present invention, this implicit serialization is made possible without any extra messages, including explicit lock messages of the resources. Instead, a member of the processing group initiates a transaction with the PTC message. It then waits  
10 for the acknowledgment that the other members received the PTC message, and this acknowledgment is called the PTC\_ACK message. When the initiating member receives all of the PTC\_ACKs, it can then issue the commit message. Therefore, concurrent transactions are serialized by making the group  
15 coordinator hold its acknowledgment, if it detects conflicts in the PTC phase.

Thus, the conflict problem depicted in the previous example is solved as follows (assuming Server 3 is the coordinator):

20	<u>Server 1</u>	<u>Server 2</u>	<u>Server 3</u>
	PTC(T1{A})	PTC(T1{A})	PTC(T2{A})
	PTC(T2{A})	PTC(T2{A})	PTC(T1{A}) *coordinator
	detects simultaneous use of token "A"		

/\*\*The servers wait for the acknowledgment that the PTCs  
25 were received

PTC\_ACK(T2{A}) \*coordinator  
acknowledges only receiving  
T2 although it has already  
received T1)

5 CMT(T2) CMT(T2) CMT(T2) \*all members commit T2  
PTC\_ACK(T1{A}) \*coordinator  
now acknowledges receiving  
T1

CMT(T1) CMT(T1) CMT(T1) \*all members commit T1

10 During the two-phase commit process (and other  
processing) of a distributed transaction, a failure may  
occur. If such a failure occurs, procedures are in place  
for recovery therefrom, in accordance with an aspect of the  
present invention. In one example, a transparent recovery  
15 of the DSTS system is performed, and no pending transactions  
are lost during the recovery process. As one example, the  
pending transactions are completed without requiring the  
reposting of the transactions, even if a number of members  
of the DSTS group fail.

20 In accordance with an aspect of the present invention,  
a facility is provided that makes possible the completion of  
a pending transaction in the event that any member of the  
DSTS group experiences a failure. Since the DSTS system can  
recover from the failure of one or more of the member copies  
25 of the system, the system is said to be highly available.  
The solution to this problem is complicated by the fact  
that, even though the DSTS system guarantees that  
transactions complete synchronously, the arrival of the

messages in a two-phase protocol is not synchronous. That is, not all the members receive the PTC and CMT messages at the same time, and as a consequence at any point in time, each member may have received a different set of messages  
5 related to a protocol, and the messages may have been received in different order.

For example, consider a snapshot of the DSTS taken during normal operation at T=4, in FIG. 25. At that point, each server has received the following set of messages:

10

<u>Server 1</u>	<u>Server 2</u>	<u>Server 3</u>
PTC (A)	PTC (B)	PTC (C)
PTC (B)	PTC (A)	PTC (A)
CMT (A)	PTC (C)	

15

Now, assume that Server 2 failed at T=4.

In the event of a failure, one of the surviving members is elected as a group coordinator. In this example, it is assumed that Server 1 is elected as the group coordinator. The group coordinator participates in recovery, as described  
20 herein.

One embodiment of the logic associated with a recovery facility is described with reference to FIG. 26. Initially, each surviving member sends to the group coordinator a list of the transaction identifiers for which PTCs were observed,  
25 since the last synchronization point, STEP 2600. In this example, Server 3 sends PTC(C) and PCT(A). Subsequently,

the group coordinator compares the PTC identifiers sent by the other surviving member(s) with its own list of PTCs, STEP 2602. In this example, the list from Server 3 is compared against {PTC(B) and PTC(A)}.

5       Next, the group coordinator requests the actual PTC message for any message that was reported by other members, but not received by the coordinator, STEP 2604. For example, the group coordinator, Server 1, requests from Server 3, PTC(C) message. At this point, the group  
10 coordinator has knowledge of all pending transactions, since the last synchronization point. The group coordinator now assumes the role of protocol initiator for all pending protocols. The other members of the group know that the protocol initiator role was changed because the system goes  
15 into recovery mode when a failure occurs.

The group coordinator sends PTC messages to any other surviving members, for all the PTC messages that are in the union of its PTC list and the other PTC list that it received in STEP 2600, STEP 2606. For example, the group  
20 coordinator sends out {PTC(A), PTC(B), PTC(C)}. The surviving group members receive the pending PTCs, and store the ones that they have not yet received, STEP 2608. For example, Server 3 stores PTC(B).

Subsequently, the surviving members send PTC\_ACK  
25 messages for each of the PTCs that were received, STEP 2610. As the PTC\_ACKS are received for the group members for each PTC, the group coordinator sends a commit (CMT) message,

STEP 2612. As the surviving members receive the commit message, they send CMT\_ACKS messages, STEP 2614. When the CMT\_ACKS messages are received for the pending transactions, the DSTS system has reached another synchronization point  
5 (i.e., no pending transactions).

Advantageously, the details of the two-phase commit process is hidden from the client application. In particular, the client application has no knowledge that there are other copies of the application involved in the  
10 commit process.

Further, advantageously, the recovery technique described above can take more than one failure. That is, it can successfully complete transactions, even if group members continue to fail, and even if the recovery is  
15 already in progress, as long as, for instance, a quorum of the group members is maintained. When a failure is noticed, the technique is restarted from the beginning. A transaction may be lost, however, if the initiator of the transaction fails before it can send out any PTC messages,  
20 or if all of a majority of the recipients of a PTC message fails after receiving the message. The recovery technique is applicable to all types of applications, even for applications that do not support rollback operations. Further, it is a useful communications protocol for shared  
25 nothing distributed systems.

In addition to the above, a failed member can rejoin the group by having the failed member detect the last

synchronization point that is observed and obtaining from the current group the delta of transactions that it needs to reach the most recent synchronization point of the DSTS system.

- 5           In one embodiment, group membership and group state are employed in the recovery of the DSTS system.

Described above are various aspects of managing replicated distributed synchronous transactions.

- Advantageously, the replication details are hidden from the  
10 client applications (e.g., no voting in two-phase commit, no participation in group protocols). One or more of the aspects of the present invention are applicable to homogeneous systems, as well as heterogeneous systems. As one example, capabilities are provided to facilitate the  
15 interoperability of the systems of a heterogeneous environment.

- The present invention can be included in an article of manufacture (e.g., one or more computer program products) having, for instance, computer usable media. The media has  
20 embodied therein, for instance, computer readable program code means for providing and facilitating the capabilities of the present invention. The article of manufacture can be included as a part of a computer system or sold separately.

- Additionally, at least one program storage device  
25 readable by a machine, tangibly embodying at least one

program of instructions executable by the machine to perform the capabilities of the present invention can be provided.

The flow diagrams depicted herein are just examples. There may be many variations to these diagrams or the steps  
5 (or operations) described therein without departing from the spirit of the invention. For instance, the steps may be performed in a differing order, or steps may be added, deleted or modified. All of these variations are considered a part of the claimed invention.

10 Although preferred embodiments have been depicted and described in detail herein, it will be apparent to those skilled in the relevant art that various modifications, additions, substitutions and the like can be made without departing from the spirit of the invention and these are  
15 therefore considered to be within the scope of the invention as defined in the following claims.



1           2.    A system of serializing replicated transactions in  
2   a distributed computing environment, said system comprising:

3                means for initiating a modification operation on a  
4   resource of a distributed computing environment;

5                means for detecting whether a conflict for said  
6   resource exists, during a phase of said modification  
7   operation; and

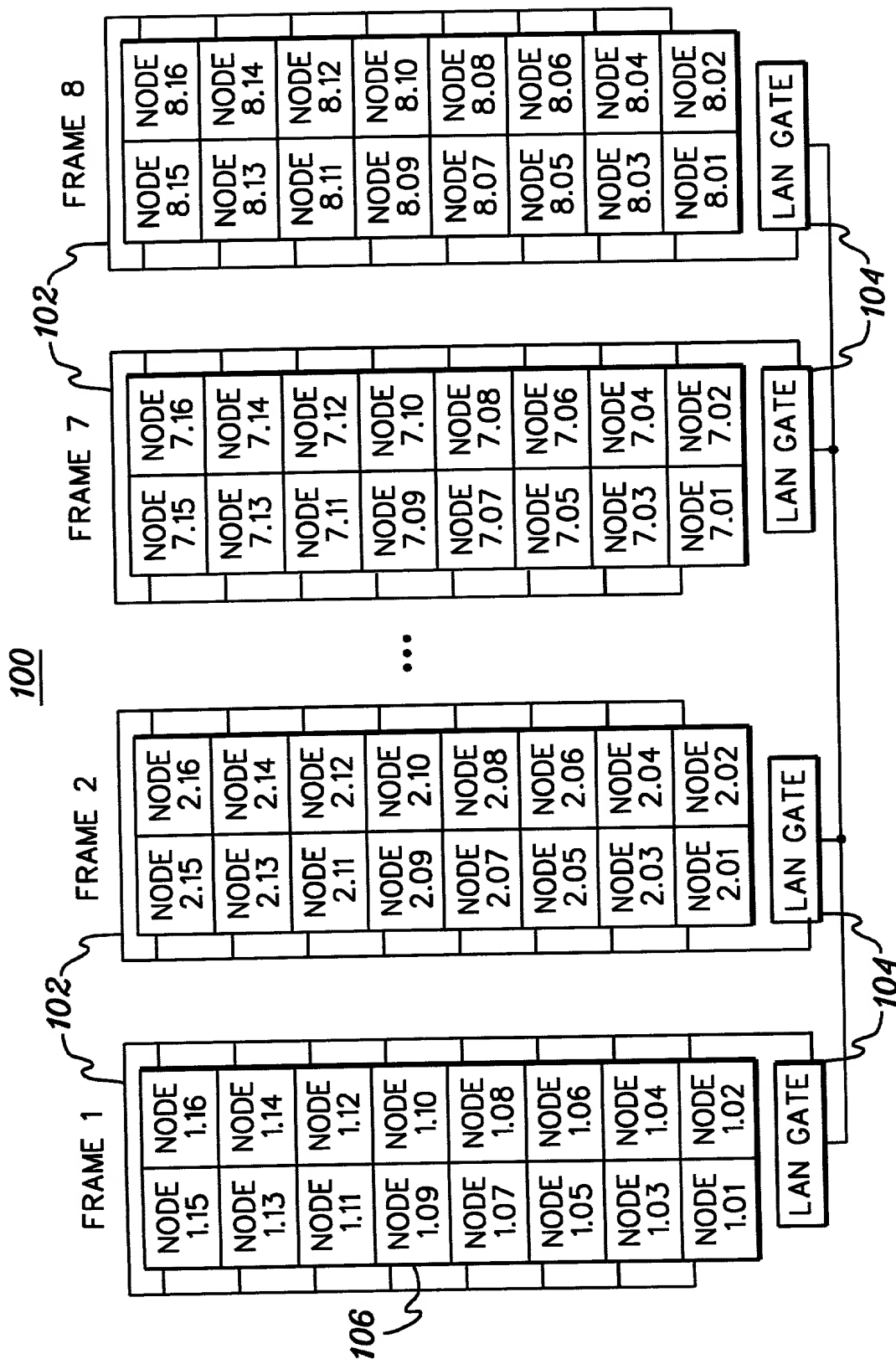
8                means for satisfying said conflict, if said  
9   conflict exists, without requiring explicit locking of  
10   said resource.

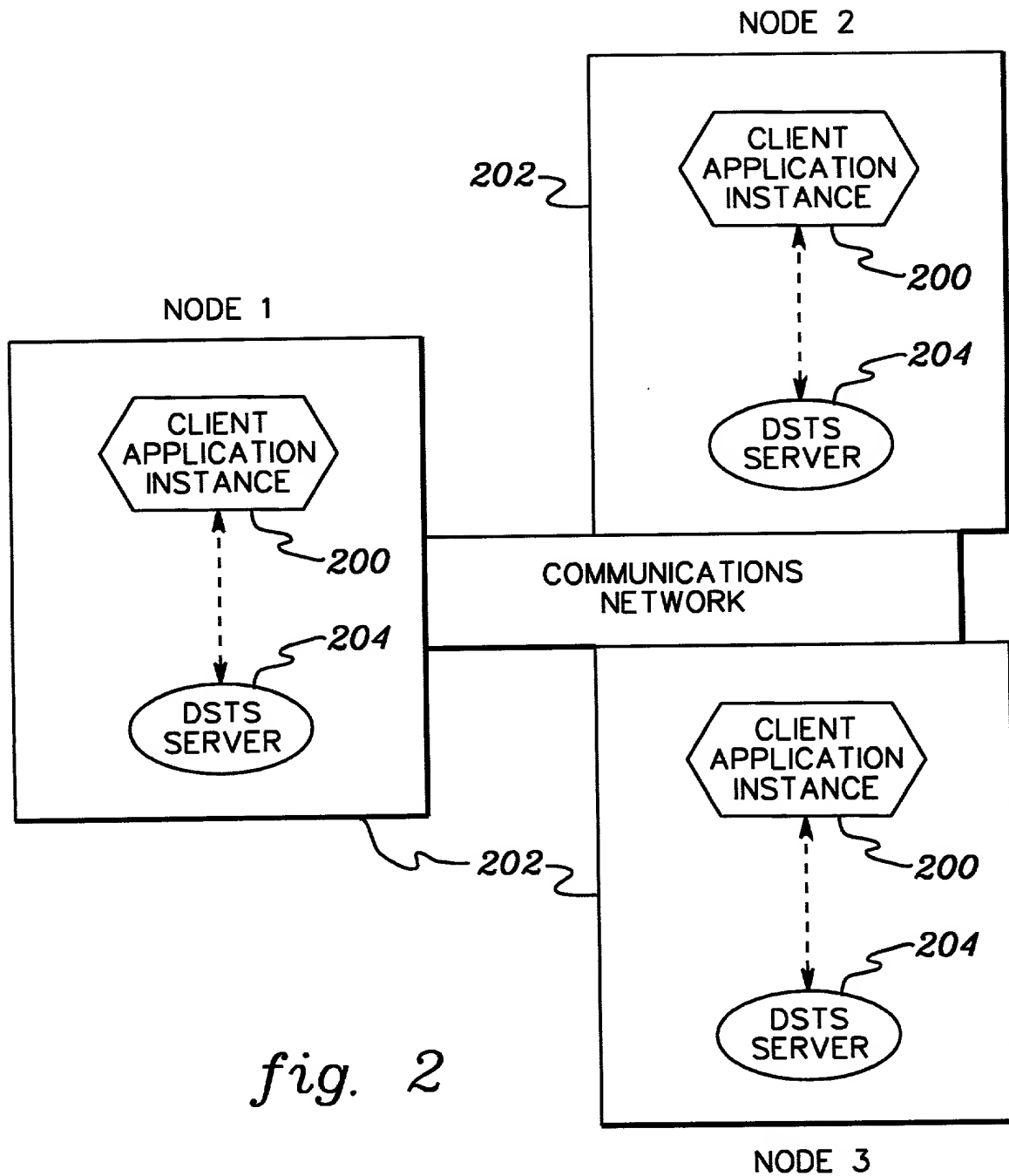


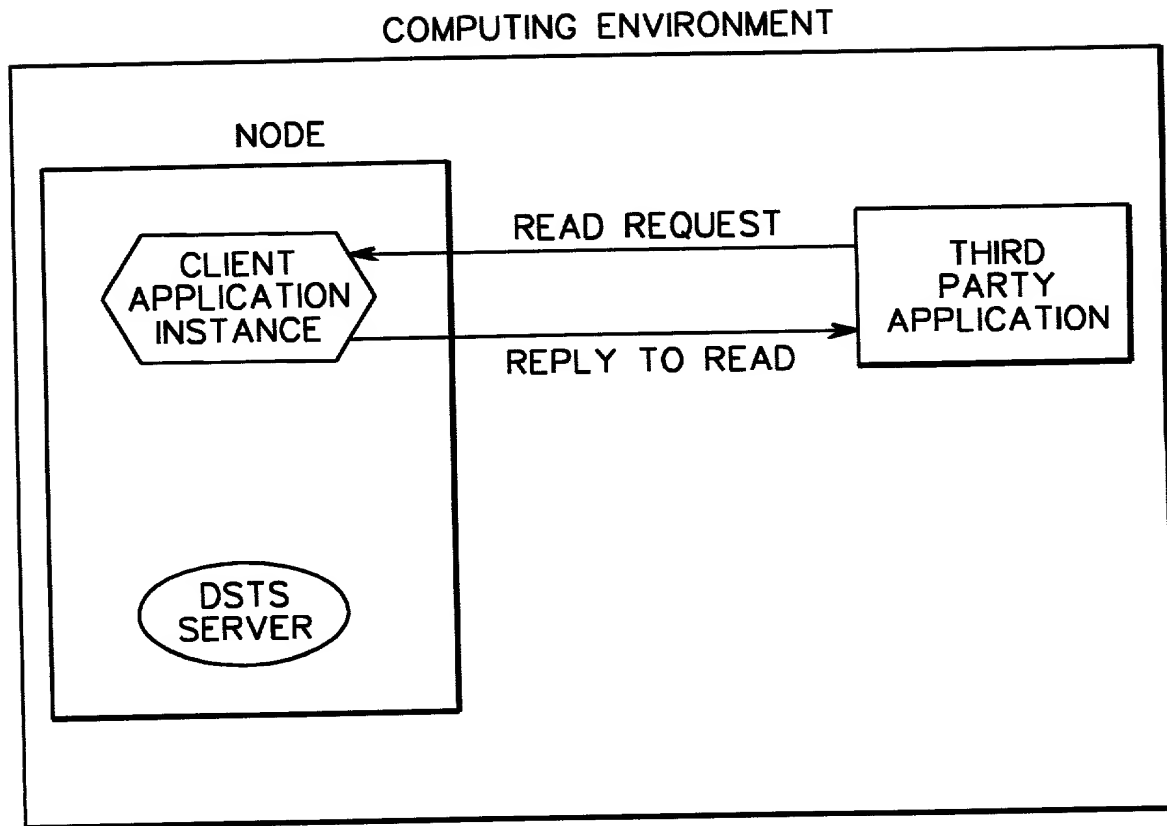
**METHOD, SYSTEM AND PROGRAM PRODUCTS FOR  
SERIALIZING REPLICATED TRANSACTIONS OF A  
DISTRIBUTED COMPUTING ENVIRONMENT**

**Abstract of the Disclosure**

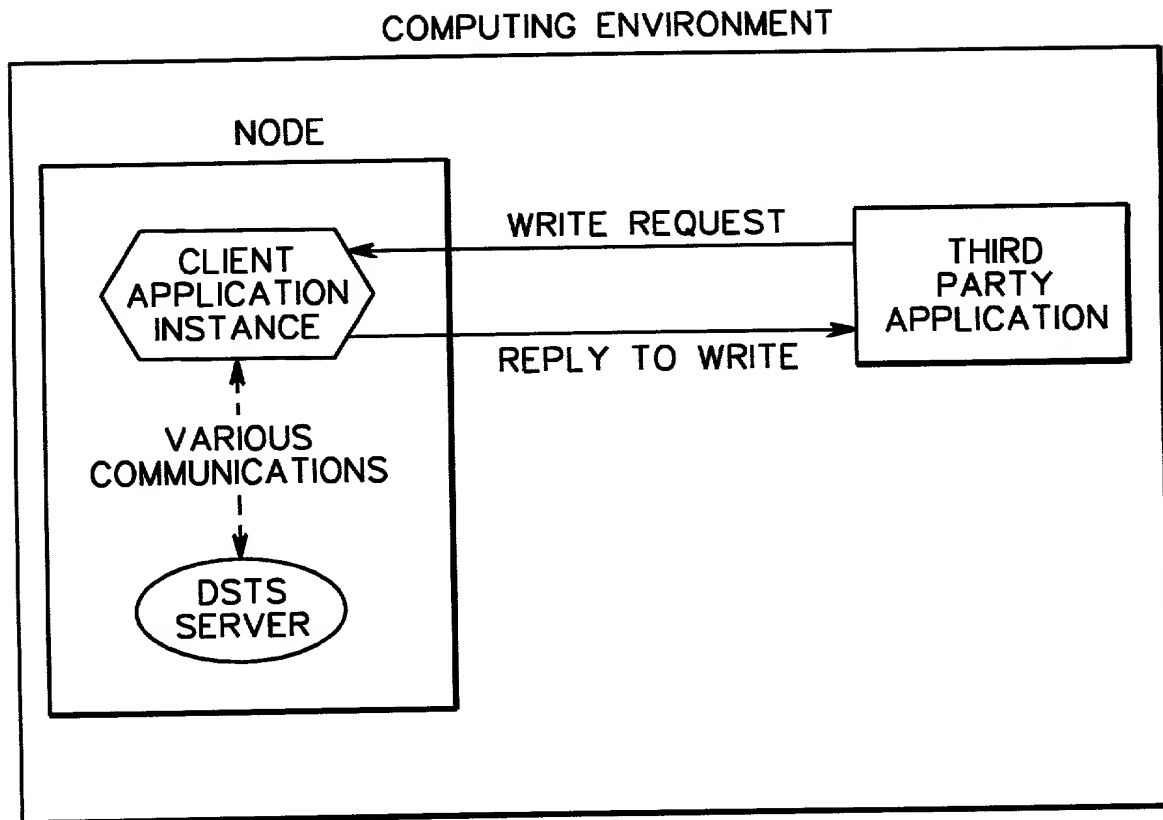
5           The management and use of replicated distributed  
transactions is facilitated. A distributed synchronous  
transaction system protocol is provided to manage the  
replication of distributed transactions for client  
application instances. The distributed synchronous  
10 transaction system allows transactions to be replicated  
without having the client application instances be aware of  
other instances to receive the transaction. Further, if a  
failure occurs during processing of a distributed replicated  
transaction, the distributed synchronous transaction system  
15 manages the recovery of the failure.



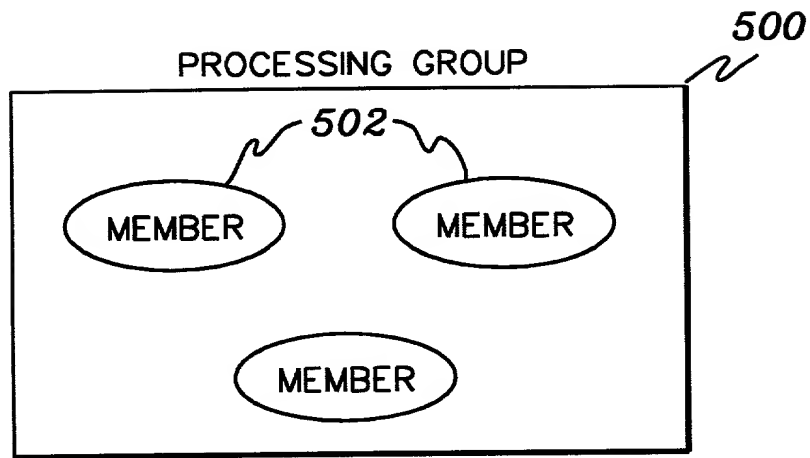




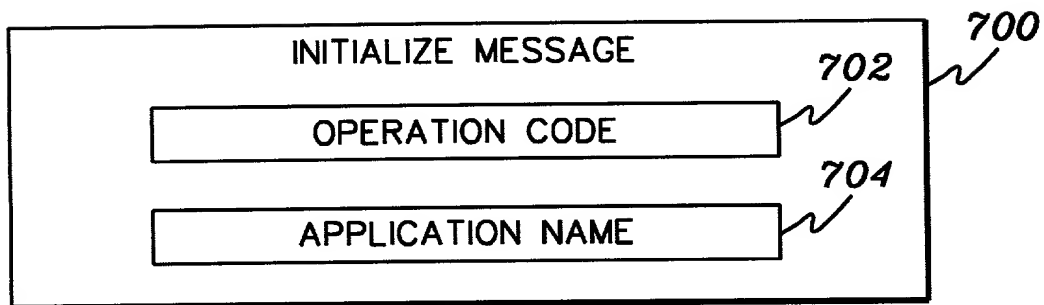
*fig. 3*



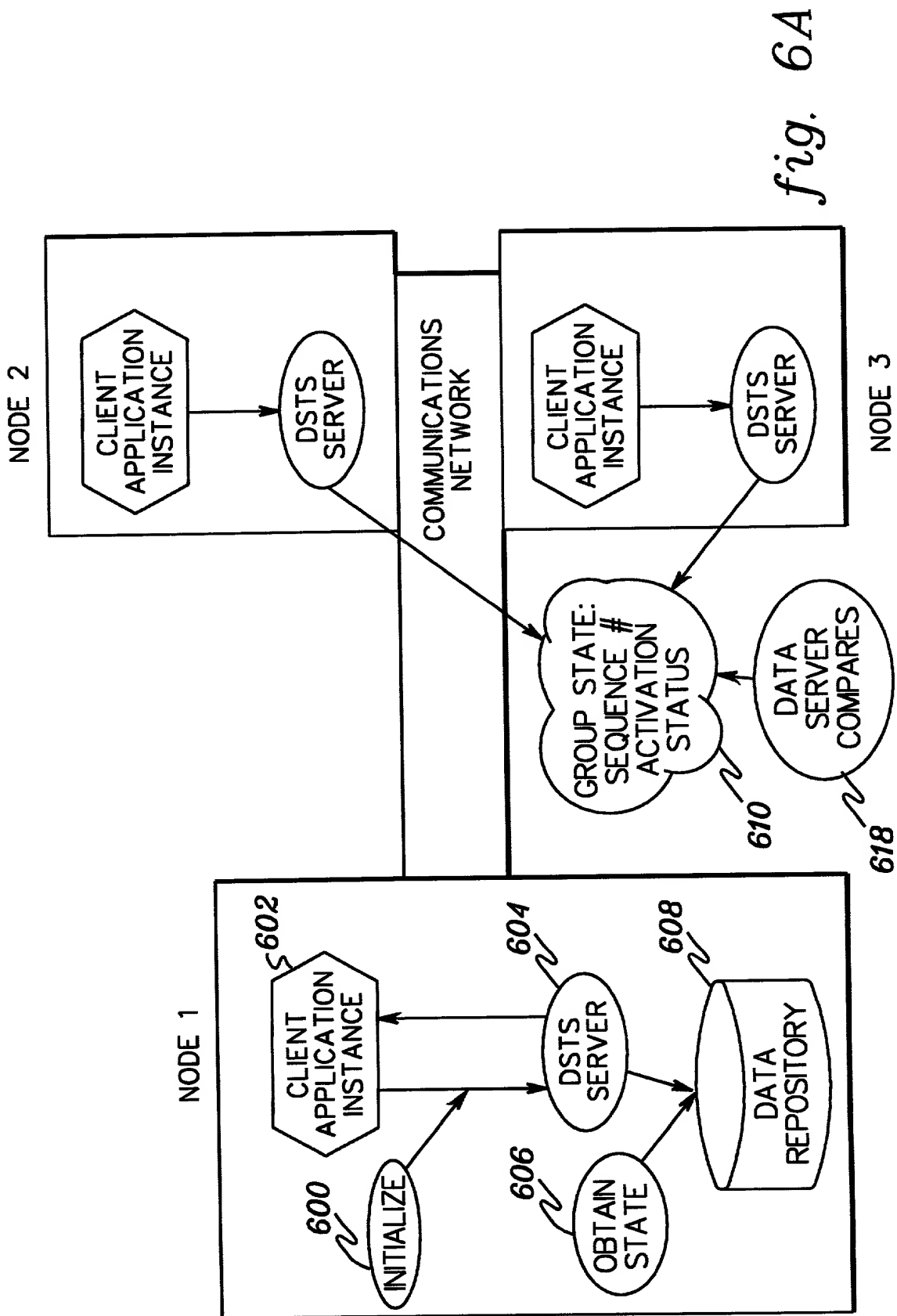
*fig. 4*



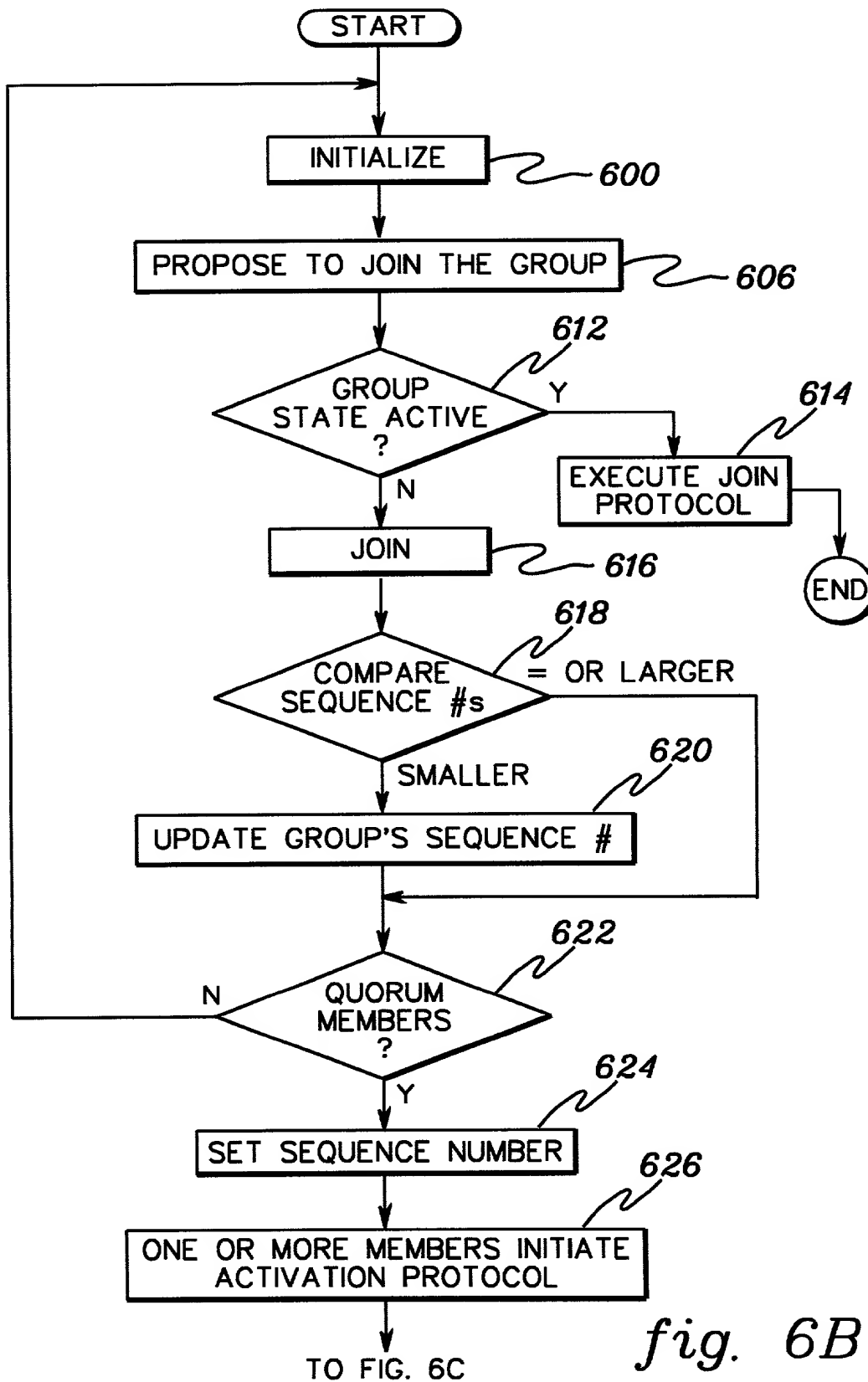
*fig. 5*



*fig. 7*



ACTIVATION - INITIALIZATION STEP



*fig. 6C*

```

graph TD
    628[RECEIVE ACTIVATION MESSAGE] --> 630[OBTAIN COPY OF GROUP STATE, IF NEEDED]
    630 --> 632[REINITIALIZE, IF NEEDED]
    632 --> 634[PROPOSE TO GO TO SECOND PHASE]
    634 --> 6D[TO FIG. 6D]

```

FROM FIG. 6C

*fig. 6D*

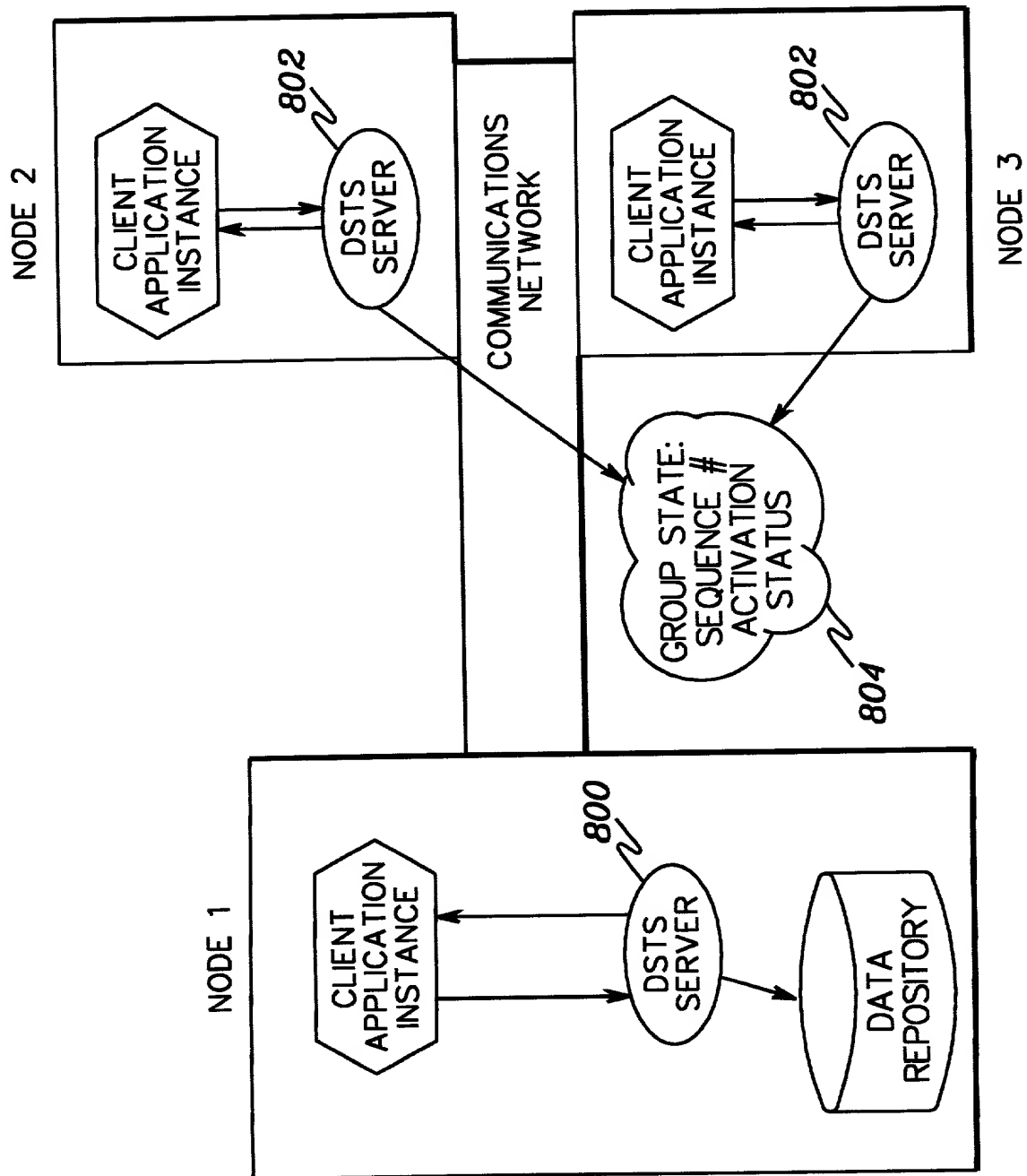


fig. 8

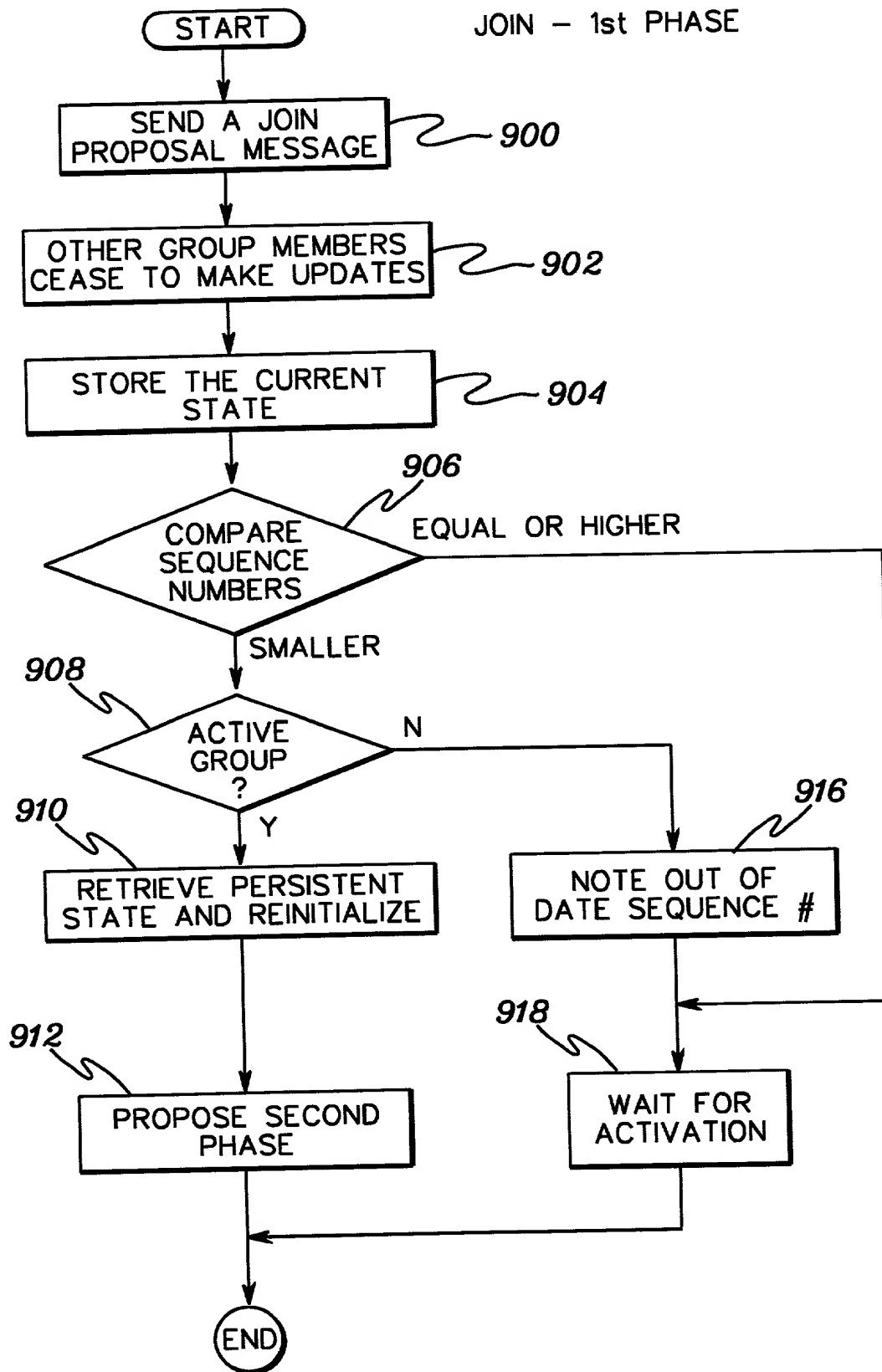
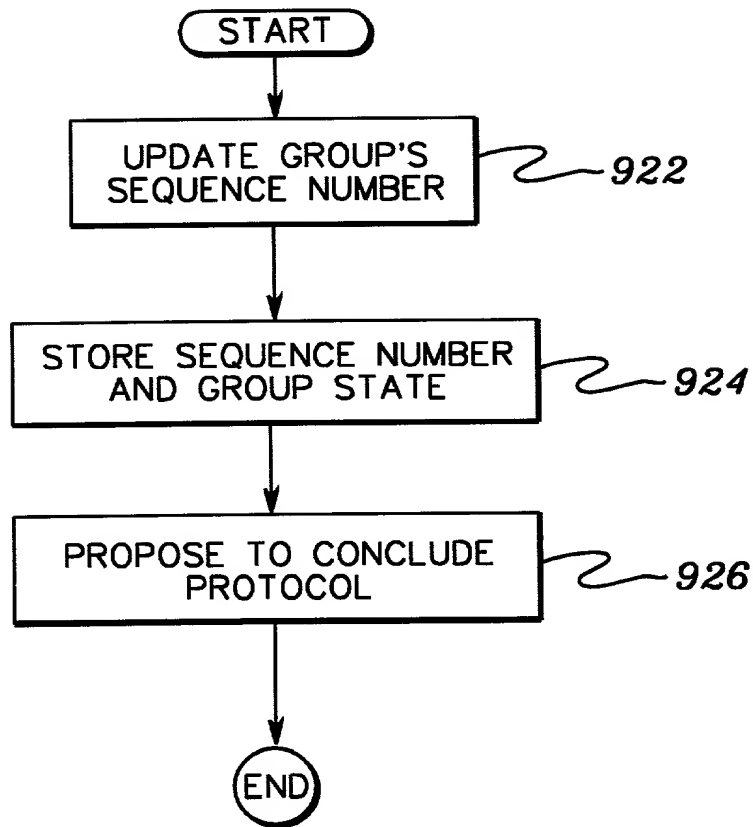
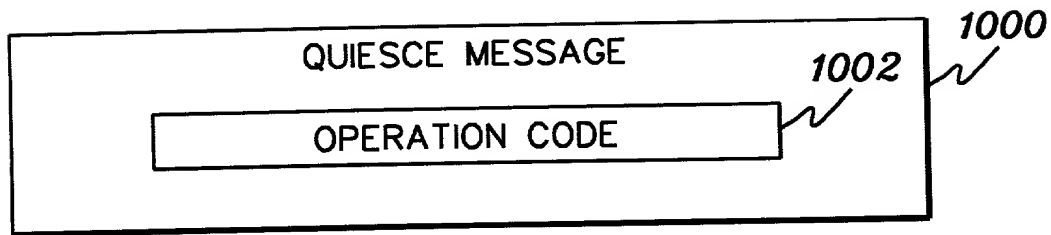


fig. 9A

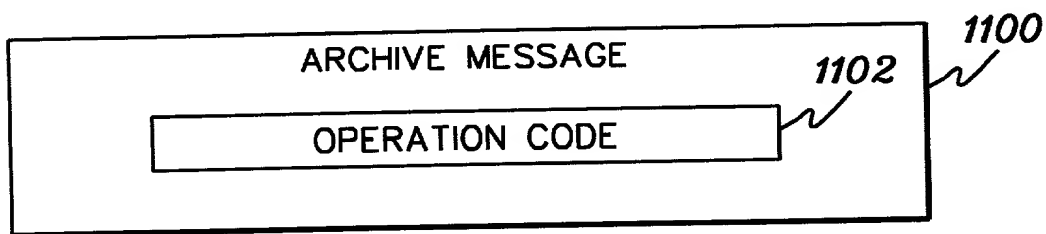
JOIN - 2nd PHASE



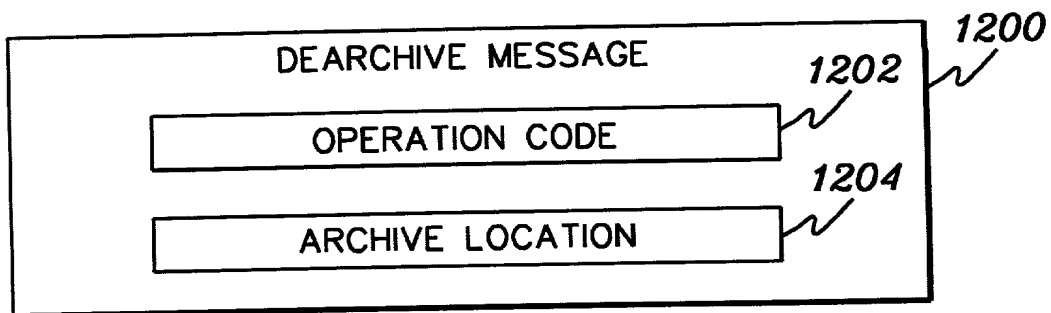
*fig. 9B*



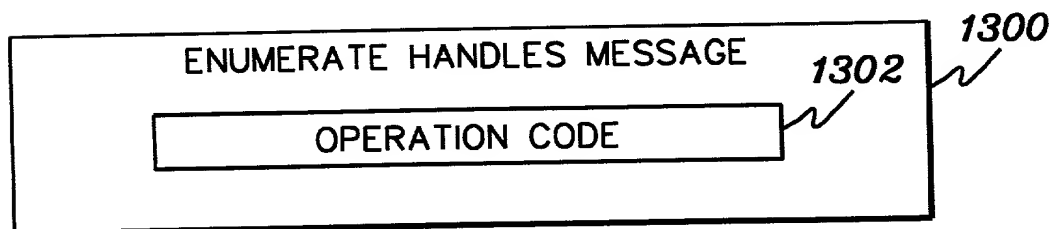
*fig. 10*



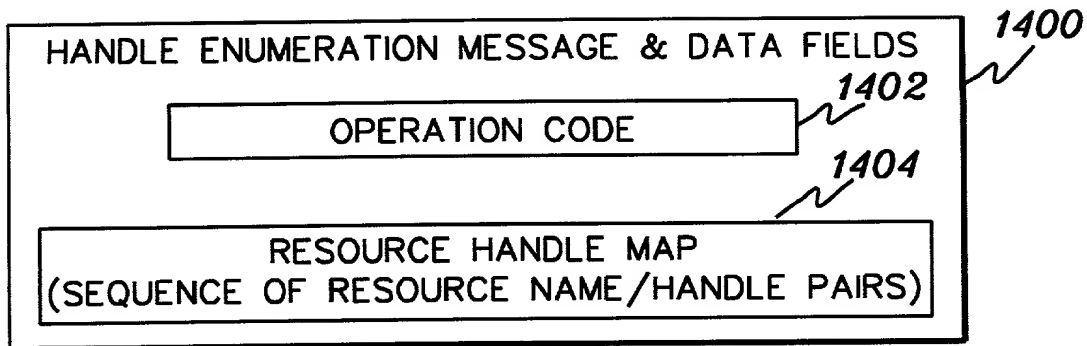
*fig. 11*



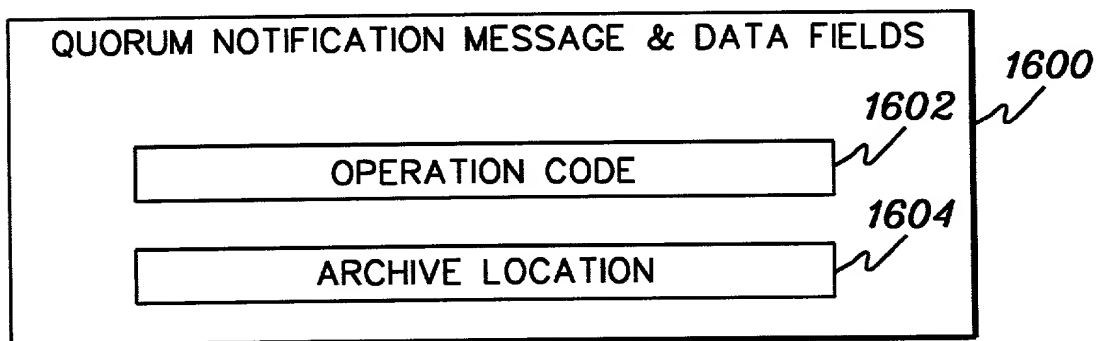
*fig. 12*



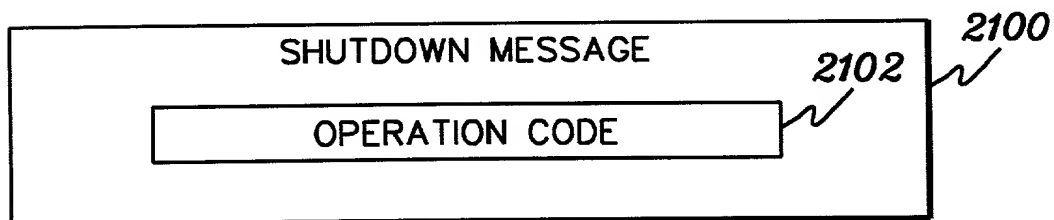
*fig. 13*



*fig. 14*



*fig. 16*



*fig. 21*

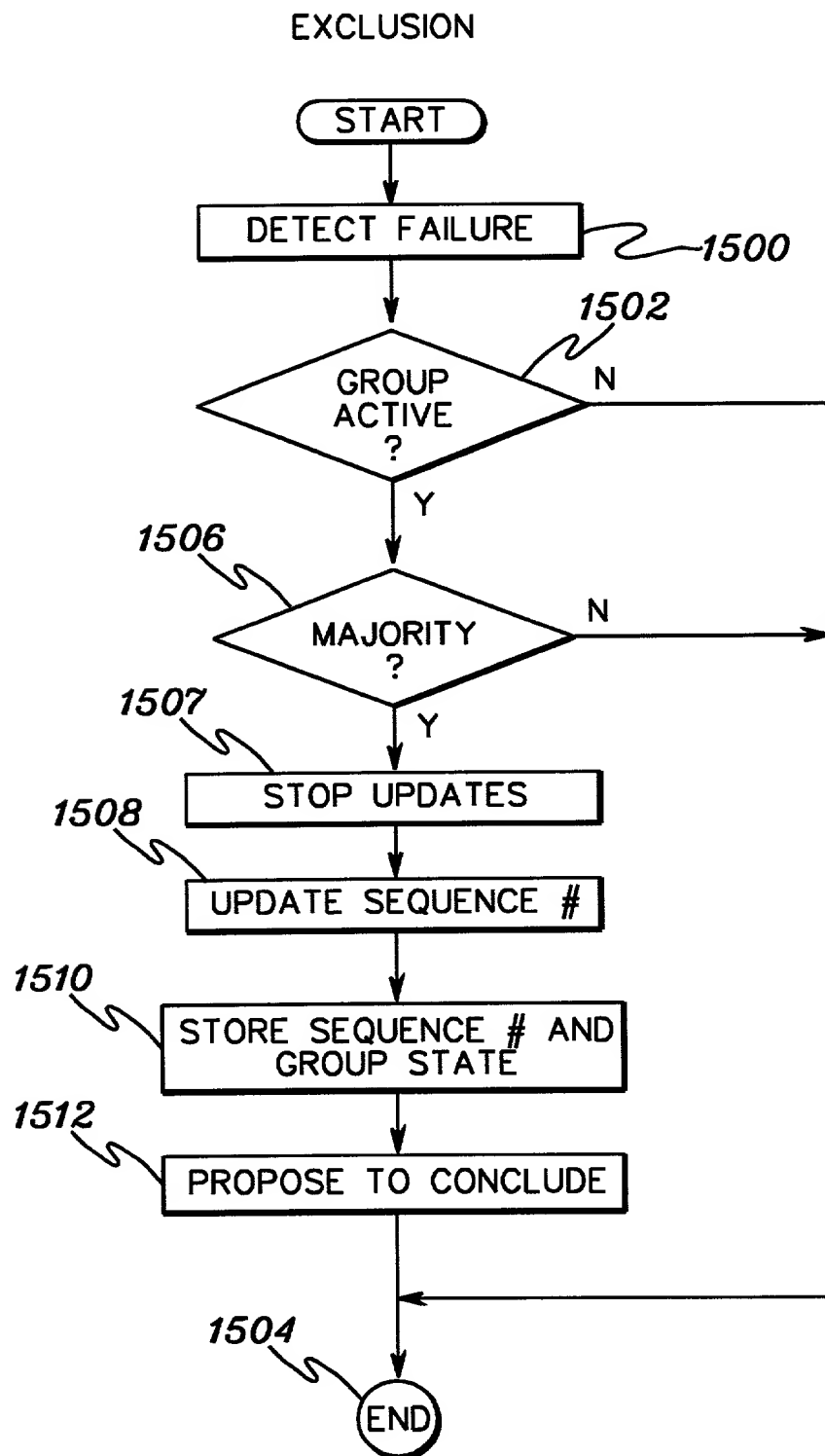
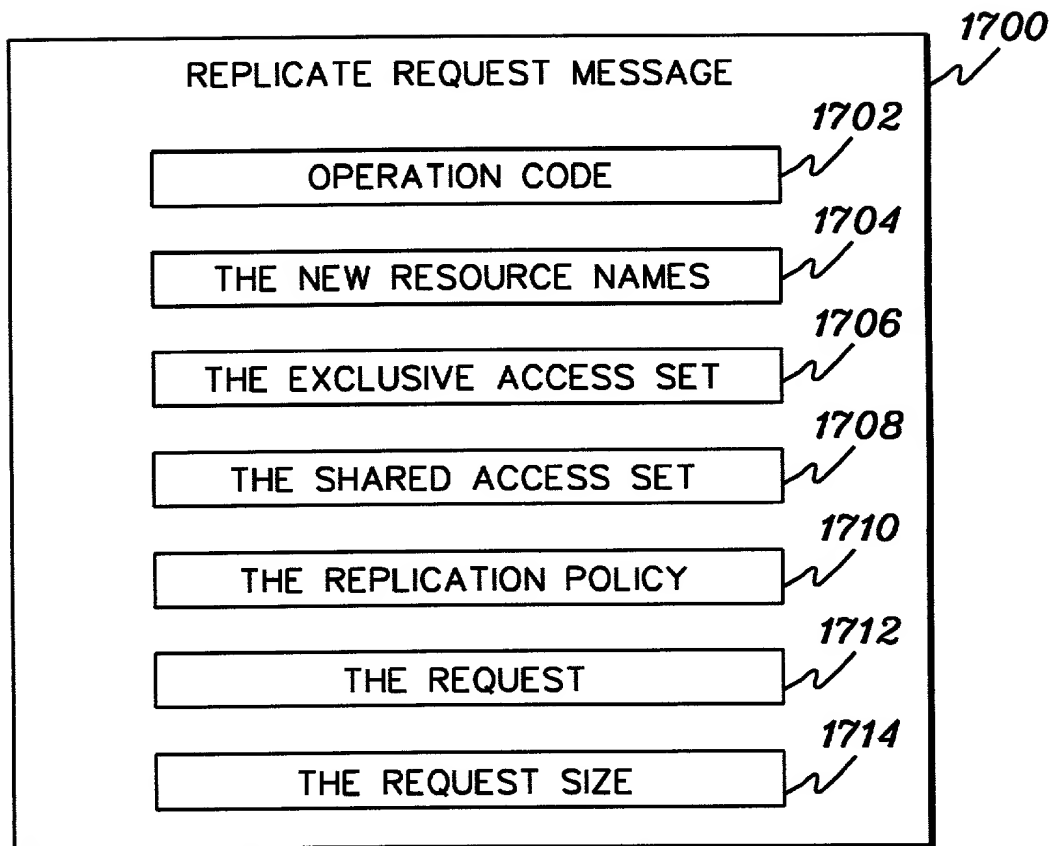
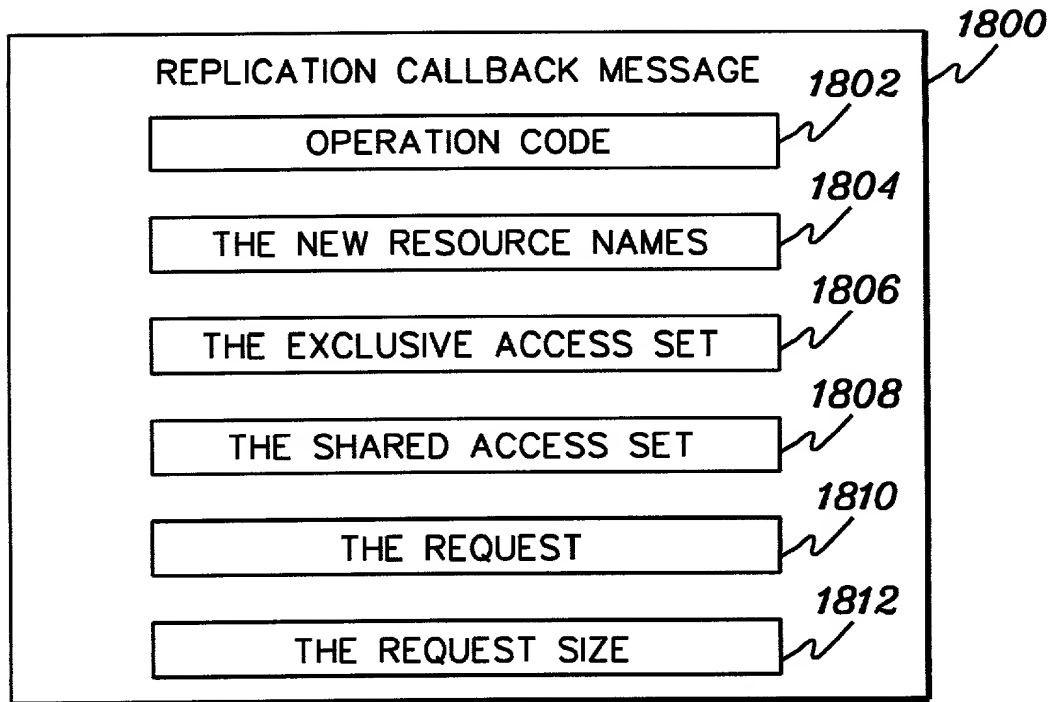


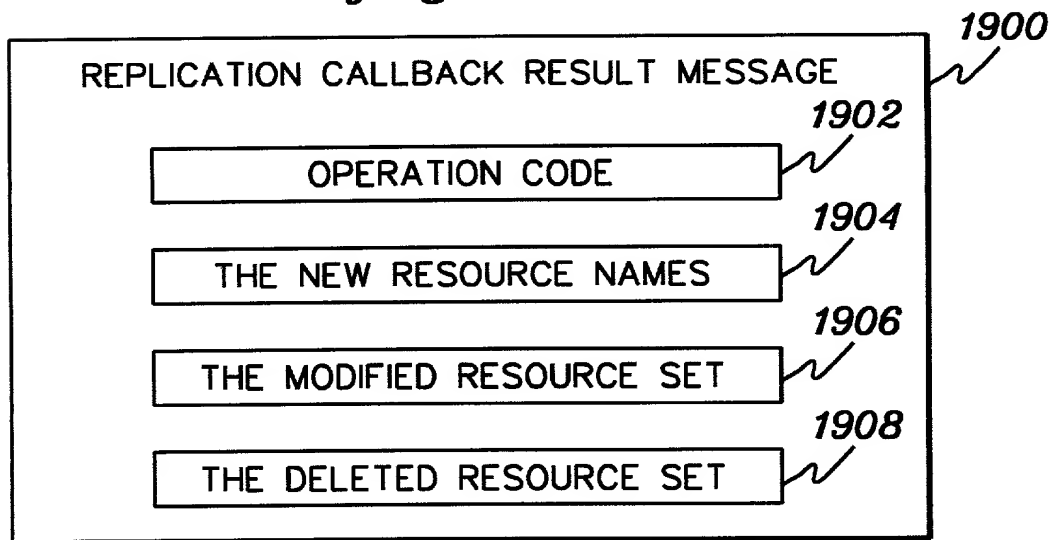
fig. 15



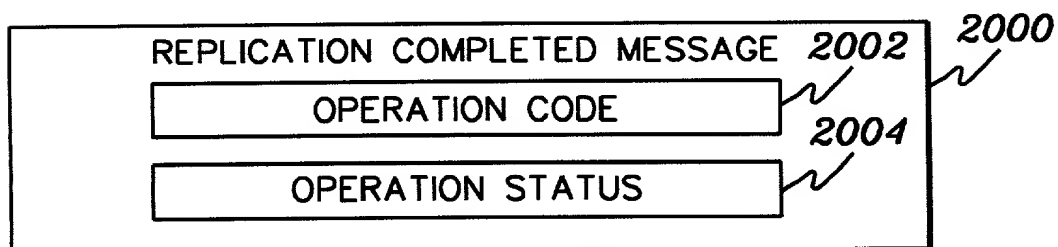
*fig. 17*



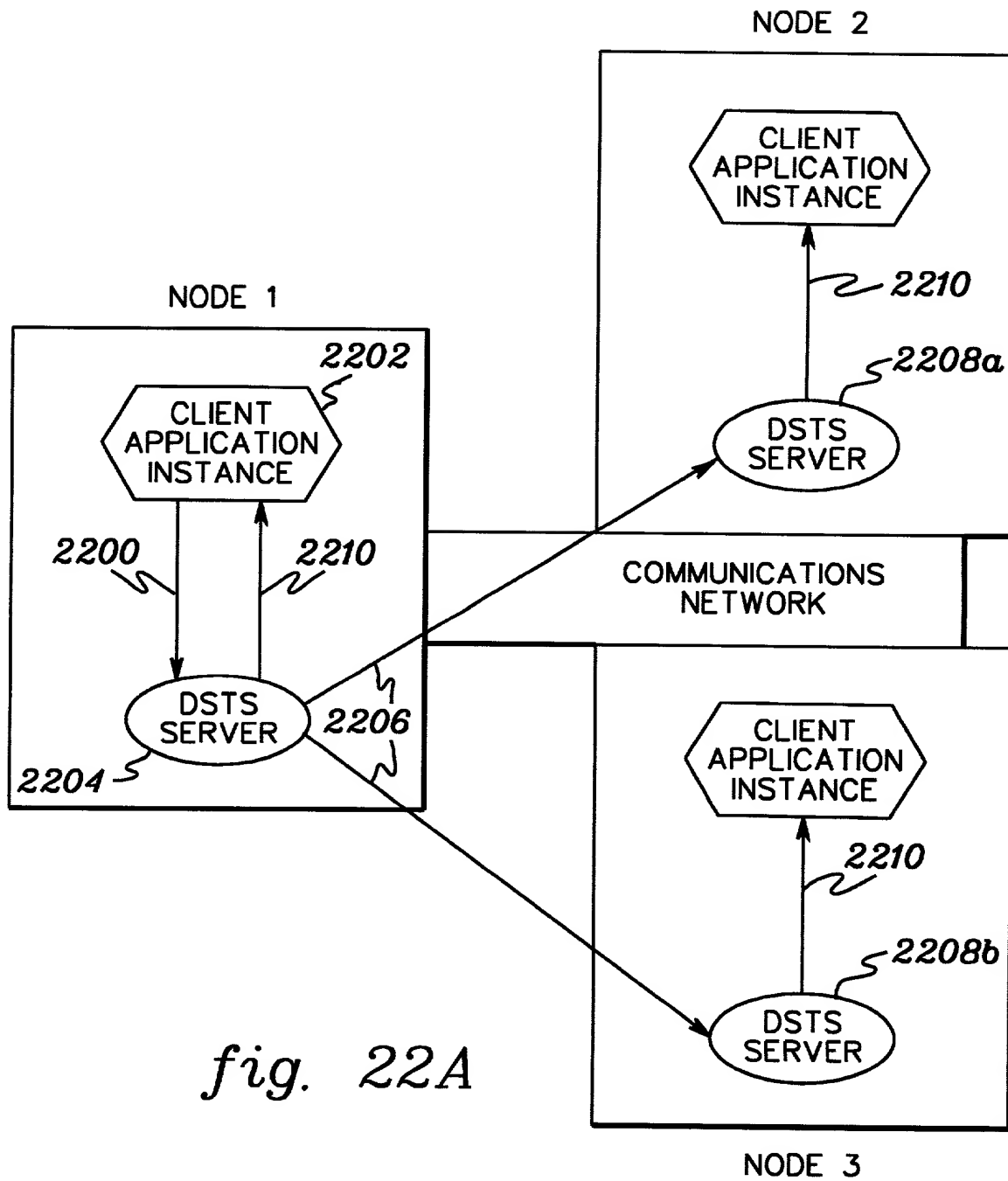
*fig. 18*

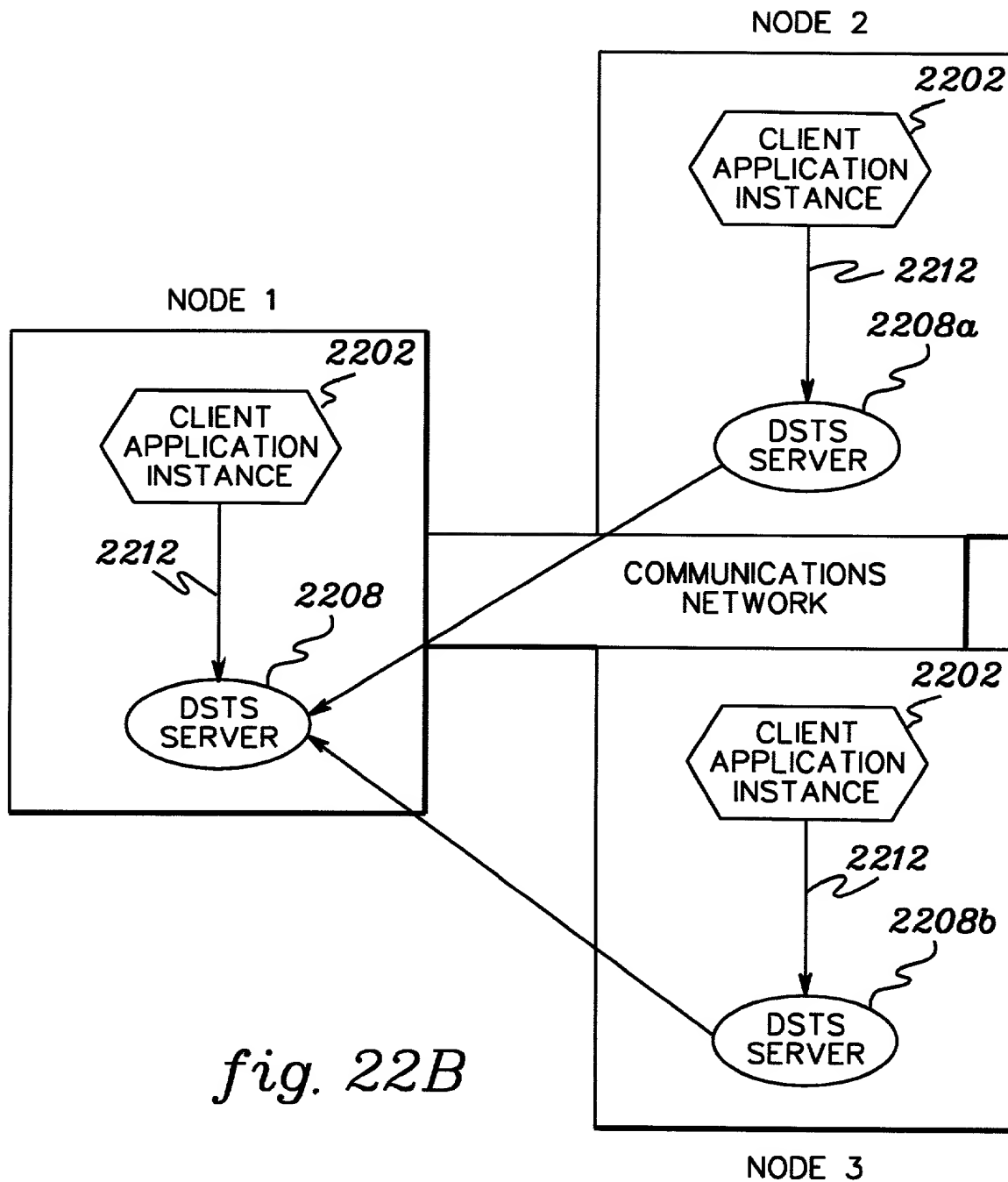


*fig. 19*

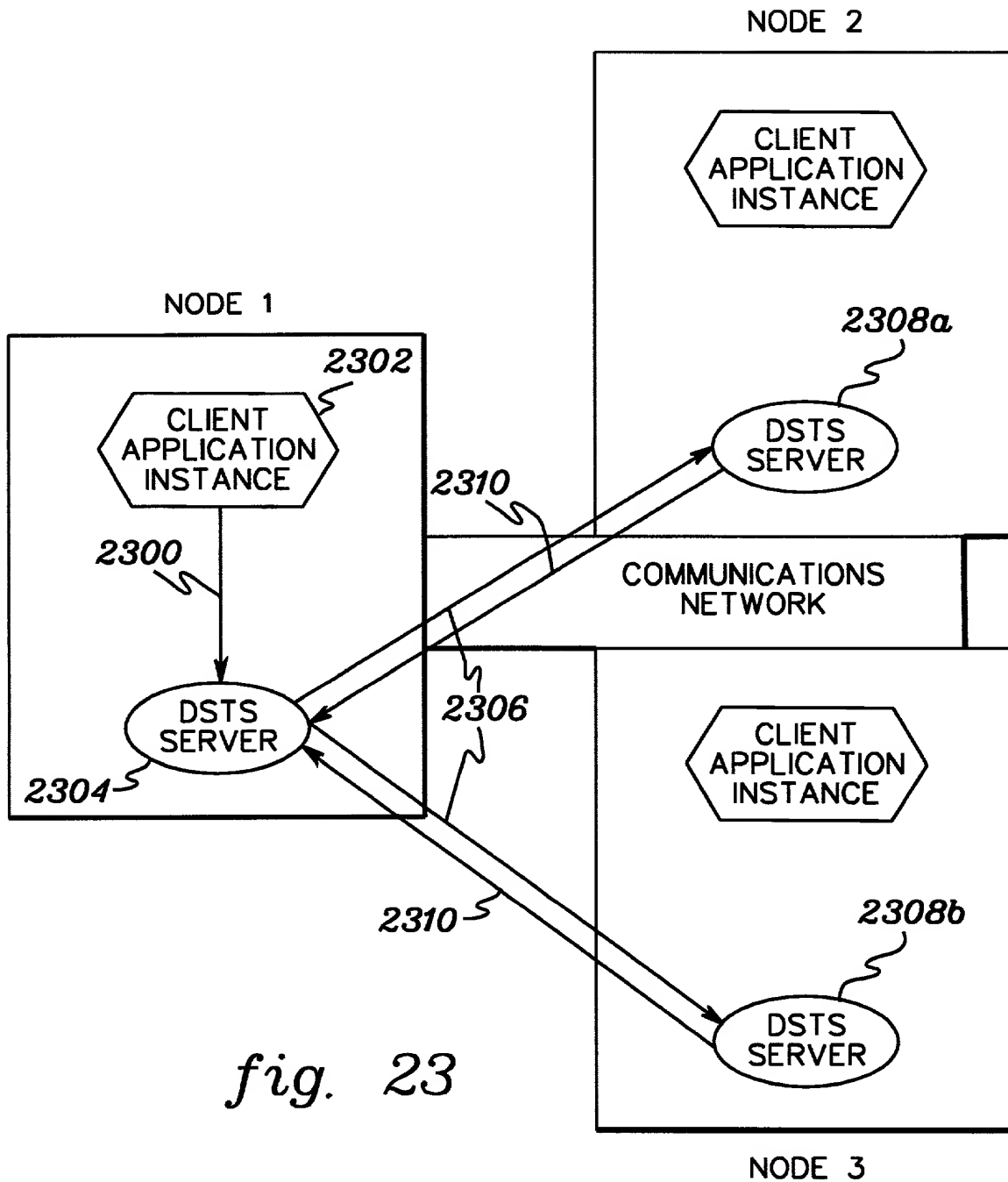


*fig. 20*

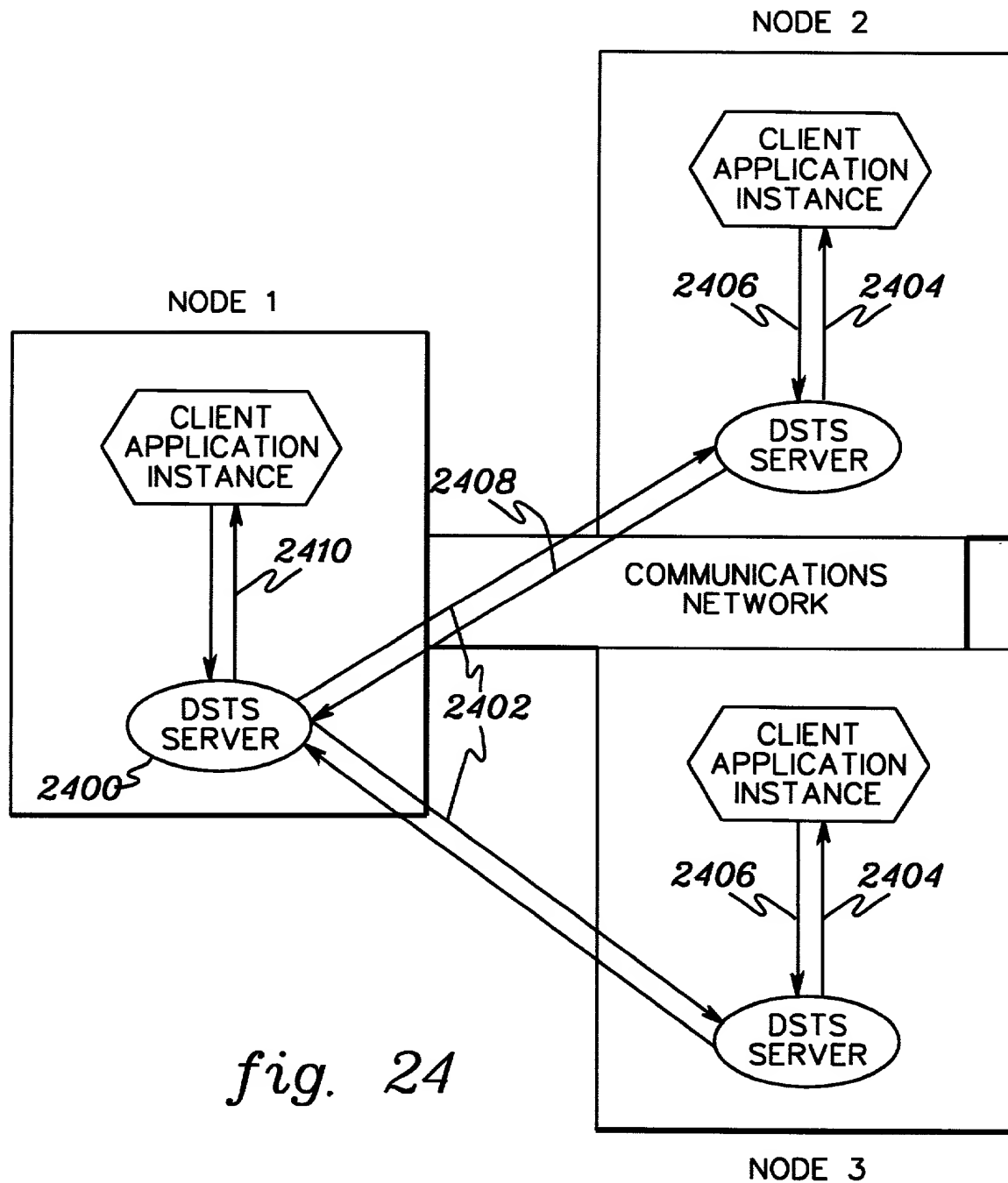


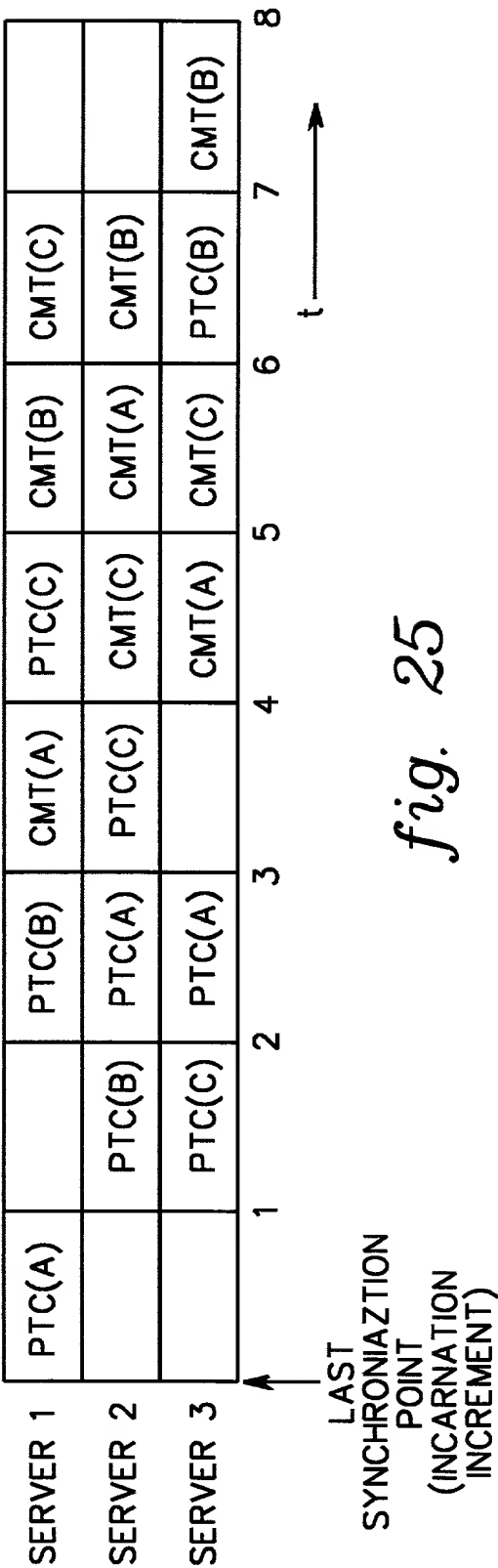


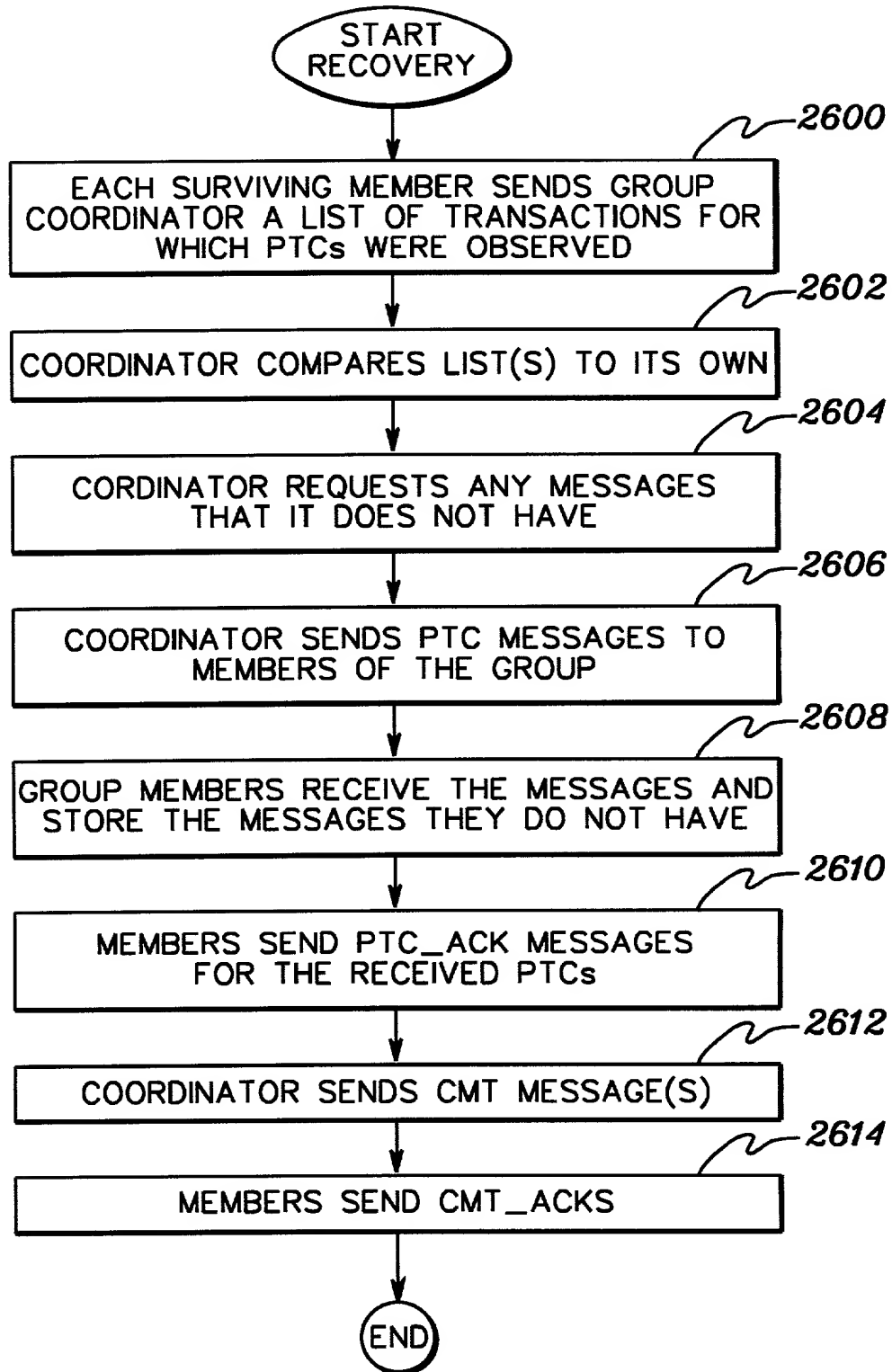
*fig. 22B*



*fig. 23*







*fig. 26*

Docket No.  
POU9-2000-0014-US1

# Declaration and Power of Attorney For Patent Application

## English Language Declaration

As a below named inventor, I hereby declare that:

My residence, post office address and citizenship are as stated below next to my name,

I believe I am the original, first and sole inventor (if only one name is listed below) or an original, first and joint inventor (if plural names are listed below) of the subject matter which is claimed and for which a patent is sought on the invention entitled

**METHOD, SYSTEM AND PROGRAM PRODUCTS FOR SERIALIZING REPLICATED  
TRANSACTIONS OF A DISTRIBUTED COMPUTING ENVIRONMENT**

the specification of which

(check one)

☒ is attached hereto.

☐ was filed on \_\_\_\_\_ as United States Application No. or PCT International  
Application Number \_\_\_\_\_  
and was amended on \_\_\_\_\_  
(if applicable)

I hereby state that I have reviewed and understand the contents of the above identified specification, including the claims, as amended by any amendment referred to above.

I acknowledge the duty to disclose to the United States Patent and Trademark Office all information known to me to be material to patentability as defined in Title 37, Code of Federal Regulations, Section 1.56.

I hereby claim foreign priority benefits under Title 35, United States Code, Section 119(a)-(d) or Section 365(b) of any foreign application(s) for patent or inventor's certificate, or Section 365(a) of any PCT International application which designated at least one country other than the United States, listed below and have also identified below, by checking the box, any foreign application for patent or inventor's certificate or PCT International application having a filing date before that of the application on which priority is claimed.

Prior Foreign Application(s)

Priority Not Claimed

(Number)

(Country)

(Day/Month/Year Filed)

☐

(Number)

(Country)

(Day/Month/Year Filed)

☐

(Number)

(Country)

(Day/Month/Year Filed)

☐

I hereby claim the benefit under 35 U.S.C. Section 119(e) of any United States provisional application(s) listed below:

\_\_\_\_\_  
(Application Serial No.)

\_\_\_\_\_  
(Filing Date)

\_\_\_\_\_  
(Application Serial No.)

\_\_\_\_\_  
(Filing Date)

\_\_\_\_\_  
(Application Serial No.)

\_\_\_\_\_  
(Filing Date)

I hereby claim the benefit under 35 U. S. C. Section 120 of any United States application(s), or Section 365(c) of any PCT International application designating the United States, listed below and, insofar as the subject matter of each of the claims of this application is not disclosed in the prior United States or PCT International application in the manner provided by the first paragraph of 35 U.S.C. Section 112, I acknowledge the duty to disclose to the United States Patent and Trademark Office all information known to me to be material to patentability as defined in Title 37, C. F. R., Section 1.56 which became available between the filing date of the prior application and the national or PCT International filing date of this application:

\_\_\_\_\_  
(Application Serial No.)

\_\_\_\_\_  
(Filing Date)

\_\_\_\_\_  
(Status)  
(patented, pending, abandoned)

\_\_\_\_\_  
(Application Serial No.)

\_\_\_\_\_  
(Filing Date)

\_\_\_\_\_  
(Status)  
(patented, pending, abandoned)

\_\_\_\_\_  
(Application Serial No.)

\_\_\_\_\_  
(Filing Date)

\_\_\_\_\_  
(Status)  
(patented, pending, abandoned)

I hereby declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true; and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code and that such willful false statements may jeopardize the validity of the application or any patent issued thereon.

POWER OF ATTORNEY: As a named inventor, I hereby appoint the following attorney(s) and/or agent(s) to prosecute this application and transact all business in the Patent and Trademark Office connected therewith. *(list name and registration number)*

Lynn L. Augspurger, Reg. No. 24,227

Lawrence D. Cutter, Reg. No. 28,501

Marc A. Ehrlich, Reg. No. 39,966

William B. Porter, Reg. No. 33,135

Floyd A. Gonzalez, Reg. No. 26,732

William A. Kinnaman, Jr., Reg. No. 27,650

Lily Neff, Reg. No. 38,254

Andrew J. Woznicki, Jr., Reg. No. 43,995

Christopher A. Hughes, Reg. No. 26,914

Edward A. Pennington, Reg. No. 32,588

John E. Hoel, Reg. No. 26,279

Joseph C. Redmond, Jr., Reg. No. 18,753

Jeff Rothenberg, Reg. No. 26,429

Kevin P. Radigan, Reg. No. 31,789

Blanche E. Schiller, Reg. No. 35,670

Send Correspondence to: **Blanche E. Schiller, Esq.**  
**HESLIN & ROTHENBERG, P.C.**  
**5 Columbia Circle**  
**Albany, NY 12203**

Direct Telephone Calls to: *(name and telephone number)*

**Blanche E. Schiller, Esq. (518) 452-5600**

Full name of sole or first inventor

**MARCOS N. NOVAES**

Sole or first inventor's signature

Date

Residence

**10 Ridge View Road, Hopewell Junction, NY 12533**

Citizenship

**Brazil**

Post Office Address

**10 Ridge View Road, Hopewell Junction, NY 12533**

Full name of second inventor, if any

**GREGORY D. LAIB**

Second inventor's signature

Date

Residence

**RD 3, 775 Oakwood Drive, Kingston, NY 12401**

Citizenship

**United States of America**

Post Office Address

**RD 3, 775 Oakwood Drive, Kingston, NY 12401**

Full name of third inventor, if any

**JEFFREY S. LUCASH**

Third inventor's signature

Date

Residence

**129 Altamont Drive, Hurley, NY 12443**

Citizenship

**United States of America**

Post Office Address

**129 Altamont Drive, Hurley, NY 12443**

Full name of fourth inventor, if any

**ROSARIO A. UCEDA-SOSA**

Fourth inventor's signature

Date

Residence

**400 High Point Drive, Apt. 301, Hartsdale, NY 10530**

Citizenship

**Spain**

Post Office Address

**400 High Point Drive, Apt. 301, Hartsdale, NY 10530**

Full name of fifth inventor, if any

Fifth inventor's signature

Date

Residence

Citizenship

Post Office Address

Full name of sixth inventor, if any

Sixth inventor's signature

Date

Residence

Citizenship

Post Office Address

OFFICE OF THE SECRETARY OF COMMERCE